

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Patent Application of: )

Koji TSUKAMOTO )

Serial No.: Not Assigned )

Filed: June 14, 2000 )

For: APPARATUS FOR RETRIEVING )  
INFORMATION USING )  
REFERENCE REASON OF )  
DOCUMENT )

Group Art Unit: Not Assigned

Examiner: Not Assigned



#2

**SUBMISSION OF CERTIFIED COPY OF PRIOR FOREIGN**  
**APPLICATION IN ACCORDANCE**  
**WITH THE REQUIREMENTS OF 37 C.F.R. §1.55**

*Honorable Commissioner of  
Patents and Trademarks  
Washington, D.C. 20231*

*Sir:*

In accordance with the provisions of 37 C.F.R. §1.55, the applicant(s) submit(s)  
herewith a certified copy of the following foreign application:

Japanese Patent Application No. 11-168552, filed: June 15, 1999.

It is respectfully requested that the applicants be given the benefit of the foreign filing  
date as evidenced by the certified papers attached hereto, in accordance with the requirements  
of 35 U.S.C. §119.

Respectfully submitted,

STAAS & HALSEY LLP

Date: June 14, 2000

By: \_\_\_\_\_

James D. Halsey, Jr.  
Registration No. 22,729

700 Eleventh Street, N.W., Suite 500  
Washington, D.C. 20001  
(202) 434-1500

PATENT OFFICE  
JAPANESE GOVERNMENT



This is to certify that the annexed is a true copy of the following application as filed with this Office.

Date of Application: June 15, 1999

Application Number: Patent Application  
No. 11-168552

Applicant(s): FUJITSU LIMITED

April 28, 2000

Commissioner,  
Patent Office

Takahiko KONDO

Certificate No.2000-3031664

8

日 本 国 特 許 庁

PATENT OFFICE  
JAPANESE GOVERNMENT

JP829 U.S. PTO  
09/594029  
06/15/00

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

1 9 9 9 年 6 月 1 5 日

出 願 番 号

Application Number:

平成 1 1 年 特 許 願 第 1 6 8 5 5 2 号

出 願 人

Applicant (s):

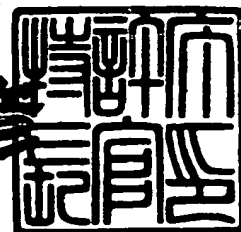
富士通株式会社

CERTIFIED COPY OF  
PRIORITY DOCUMENT

2 0 0 0 年 4 月 2 8 日

特 許 庁 長 官  
Commissioner,  
Patent Office

近 藤 隆 彦



出 証 番 号 出 証 特 2 0 0 0 - 3 0 3 1 6 6 4

【書類名】 特許願

【整理番号】 9900604

【提出日】 平成11年 6月15日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/27  
G06F 17/30

【発明の名称】 文書の参照理由を用いて情報検索を行う装置

【請求項の数】 15

【発明者】

    【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

    【氏名】 塚本 浩司

【特許出願人】

    【識別番号】 000005223

    【氏名又は名称】 富士通株式会社

【代理人】

    【識別番号】 100074099

    【住所又は居所】 東京都千代田区二番町8番地20 二番町ビル3F

    【弁理士】

    【氏名又は名称】 大菅 義之

    【電話番号】 03-3238-0031

【選任した代理人】

    【識別番号】 100067987

    【住所又は居所】 神奈川県横浜市港北区太尾町1418-305 (大倉山二番館)

    【弁理士】

    【氏名又は名称】 久木元 彰

    【電話番号】 045-545-9280

【手数料の表示】

【予納台帳番号】 012542

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9705047

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書の参照理由を用いて情報検索を行う装置

【特許請求の範囲】

【請求項 1】 与えられた文書データから被参照文書に関する文書情報を抽出する被参照文書抽出手段と、

前記文書データ内で前記被参照文書を参照している位置に関する情報を抽出する参照位置抽出手段と、

前記参照位置抽出手段により抽出された情報を解析して、前記被参照文書が参照されている理由を判断する判断手段と、

前記被参照文書抽出手段により抽出された情報と前記被参照文書が参照されている理由とを含む出力情報を出力する出力手段と

を備えることを特徴とする参照理由同定装置。

【請求項 2】 前記文書データから該文書データに関する文書情報を抽出する文書情報抽出手段をさらに備え、前記出力手段は、該文書情報抽出手段により抽出された情報を前記出力情報とともに出力することを特徴とする請求項 1 記載の参照理由同定装置。

【請求項 3】 前記判断手段は、前記被参照文書を参照している位置が属する章の情報と、該被参照文書を参照している位置の周辺の文字列の情報のうち少なくとも一方に基づいて、前記被参照文書が参照されている理由を判断することを特徴とする請求項 1 記載の参照理由同定装置。

【請求項 4】 与えられた文書データから被参照文書に関する文書情報を抽出する被参照文書抽出手段と、

前記文書データ内で前記被参照文書を参照している位置に関する情報を抽出する参照位置抽出手段と、

前記参照位置抽出手段により抽出された情報を解析して、前記被参照文書が参照されている理由を判断する判断手段と、

前記被参照文書を参照している位置の周辺の情報から、前記被参照文書のためのキーワード情報を抽出するキーワード抽出手段と、

前記被参照文書抽出手段により抽出された情報と前記被参照文書が参照されて

いる理由と前記キーワード情報とを含む出力情報を出力する出力手段と  
を備えることを特徴とするキーワード抽出装置。

【請求項 5】 文書データを格納する文書データベース手段と、  
ある文書が被参照文書を参照している理由を表す参照理由を含む参照関係を格納  
する参照関係格納手段と、

前記参照関係格納手段に格納された参照関係を利用して、前記文書データベー  
ス手段に格納された文書データの検索を行う検索手段と、

前記参照理由を含む検索結果を出力する出力手段と  
を備えることを特徴とする情報検索装置。

【請求項 6】 前記ある文書の文書データから前記被参照文書に関する文書  
情報を抽出する被参照文書抽出手段と、

前記ある文書の文書データ内で前記被参照文書を参照している位置に関する情  
報を抽出する参照位置抽出手段と、

前記参照位置抽出手段により抽出された情報を解析して、前記参照理由を判断  
する判断手段と、

前記被参照文書抽出手段により抽出された情報と前記参照理由とを含む情報を  
前記参照関係格納手段に格納する手段と

をさらに備えることを特徴とする請求項 5 記載の情報検索装置。

【請求項 7】 前記参照関係格納手段は、前記ある文書から抽出された前記  
被参照文書のためのキーワード情報を格納し、前記検索手段は、該キーワード情  
報を利用して検索を行うことを特徴とする請求項 5 記載の情報検索装置。

【請求項 8】 前記出力手段は、前記参照関係を前記参照理由に基づいてグ  
ラフィカルに表示する表示手段を含むことを特徴とする請求項 5 記載の情報検索  
装置。

【請求項 9】 前記出力手段は、前記参照関係を時系列に表示する表示手段  
を含むことを特徴とする請求項 5 記載の情報検索装置。

【請求項 10】 与えられた文書データから被参照文書に関する文書情報を  
抽出する被参照文書抽出手段と、

前記文書データ内で前記被参照文書を参照している位置に関する情報を抽出す

る参照位置抽出手段と、

前記参照位置抽出手段により抽出された情報を解析して、前記被参照文書が参照されている理由を判断する判断手段と、

複数の文書データ間において、被参照文書が参照されている理由を含む参照関係の類似度を計算し、該複数の文書データを分類する類似度判定手段と、

分類結果を出力する出力手段と

を備えることを特徴とする文書分類装置。

【請求項 11】 前記被参照文書を参照している位置の周辺の情報から、前記被参照文書のためのキーワード情報を抽出するキーワード抽出手段をさらに備え、前記類似度判定手段は、該キーワード情報を用いて前記複数の文書データを分類することを特徴とする請求項 10 記載の文書分類装置。

【請求項 12】 検索対象の情報を格納するデータベース手段と、ある情報が被参照情報を参照している理由を表す参照理由を含む参照関係を格納する参照関係格納手段と、

前記参照関係格納手段に格納された参照関係を利用して、前記データベース手段に格納された情報の検索を行う検索手段と、

前記参照理由を含む検索結果を出力する出力手段と  
を備えることを特徴とする情報検索装置。

【請求項 13】 コンピュータのためのプログラムを記録した記録媒体であって、

与えられた文書データから被参照文書に関する文書情報を抽出するステップと

前記文書データ内で前記被参照文書を参照している位置に関する情報を抽出するステップと、

前記被参照文書を参照している位置に関する情報を解析して、前記被参照文書が参照されている理由を判断するステップと、

前記被参照文書に関する文書情報と前記被参照文書が参照されている理由とを含む出力情報を出力するステップと

を含む処理を前記コンピュータに実行させるためのプログラムを記録したコン



ピュータ読み取り可能な記録媒体。

【請求項 1 4】 コンピュータのためのプログラムを記録した記録媒体であって、

ある文書が被参照文書を参照している理由を表す参照理由を含む参照関係を利用して、文書データベースに格納された文書データの検索を行うステップと、

前記参照理由を含む検索結果を出力するステップと  
を含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項 1 5】 コンピュータのためのプログラムを記録した記録媒体であって、

与えられた文書データから被参照文書に関する文書情報を抽出するステップと

前記文書データ内で前記被参照文書を参照している位置に関する情報を抽出するステップと、

前記被参照文書を参照している位置に関する情報を解析して、前記被参照文書が参照されている理由を判断する判断ステップと、

複数の文書データ間において、被参照文書が参照されている理由を含む参照関係の類似度を計算し、該複数の文書データを分類するステップと、

分類結果を出力するステップと  
を含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0 0 0 1】

【産業上の利用分野】

本発明は、文書等の情報を検索する情報検索に係り、特に、文書間に参照／被参照関係が存在する場合に、参照理由を同定し、参照理由を用いて文書を検索する装置に関する。

【0 0 0 2】

【従来の技術】

文書間に参照／被参照関係が存在する場合に、参照関係を検索に利用する技術がいくつか提案されている。このような技術としては、以下のような公開特許公報が挙げられる。

(1) 特開昭 6 3 - 2 2 8 2 2 1 (三菱電機)

文献の参照関係を記憶しておき、手がかりとなる何らかの情報からファジィ論理演算を用いて検索を行う。

(2) 特開昭 6 3 - 1 5 3 6 3 0 (日本電気)

文献の被引用関係を用いて、引用文献を検索項目とすることにより、共引用関係の文書の検索を行うことができる。ここで、共引用関係の文書とは、共通の引用文献を持つことで重要な関係にあると考えられる文献を意味する。

(3) 特開平 1 - 1 9 1 2 5 8 (リコー)

文献の本文と、本文中から自動的に抽出した参考文献名とを同時に提示して、編集を容易にする。

(4) 特開平 6 - 2 8 2 5 3 4 (日本電気)

文書が引用されたことを、引用された文書の利用者に自動的に通知する。

(5) 特開平 7 - 3 1 1 7 8 0 (キャノン)

引用関係に基づき、ある文献に関連する文献を検索する。検索結果は、重要度順に表示される。重要度は、引用回数に基づいて決められる。

(6) 特開平 8 - 2 7 2 8 1 8 (新日本製鐵)

ある文書を指定すると、それに関連する文書が表示され、表示された文書を選択することにより、さらに検索が可能になる。

(7) 特開平 9 - 1 4 6 9 6 8 (日立製作所)

ある文献が行っている参照と似たような参照を行っている他の文献を検索する。

(8) 特開平 1 0 - 1 0 5 5 7 2 (日本電気)

参照関係とキーワードをもとに、文書間に関連性があるかないかを判断し、それに基づいて文書集合を作成する。

【 0 0 0 3 】

図 2 7 は、このような従来の情報検索システムを示している。このシステムは

、検索装置 1、全文データベース 2、および参照関係データベース 3 を備える。そして、ある文献が入力されると、それと共に引用関係にある文献や似たような参照を行っている文献を、関連文献として表示する。検索装置 1 の検索部 5 は、全文データベース 2 の文献を検索し、選択部 6 は、参照関係データベース 3 の参照関係を用いて関連文献を選択する。

【0004】

関連文献の情報は、例えば、図 28 のような形式で表示される。関連文献が複数ある場合、例えば、文献の重要度の順に並べて表示される。

上述の公開特許公報のうち、特開平 8 - 2 7 2 8 1 8 と特開平 1 0 - 1 0 5 5 7 2 以外においては、基本的に、参照しているかいないかのどちらかという 2 値的な参照関係のみが用いられている。

【0005】

また、特開平 8 - 2 7 2 8 1 8 では、表示の際にどのような位置付けの参照関係なのかを表示するようにしているが、そのためには人間があらかじめ情報を付与しておく必要がある。また、特開平 1 0 - 1 0 5 5 7 2 では、参照関係とキーワードを用いて、文書間の関連性があるかないかを判断している。

【0006】

このように、従来の技術では、関連があるかないかのいずれか一方を表す 2 値情報が、文書間の関連性として用いられている。

【0007】

【発明が解決しようとする課題】

しかしながら、上述した従来の情報検索システムには、次のような問題がある。

【0008】

特開平 7 - 3 1 1 7 8 0 では、引用回数の多い文献から順に重要文献として表示しているが、利用者は装置側が一意に提示した価値基準に従って文献の重要度を逐次判断せざるを得ない。

【0009】

また、他の従来技術も、文書を参照しているかいないかの 2 値情報で参照関係

を表しており、ある文書にとって「重要であるかないか」という 1 つの観点から分析を行い、文書を検索するものである。このような 2 値情報を用いて検索を行うシステムは、次のような欠点を持っている。

- (1) 利用者になぜ重要であるのかを示すことができない。
- (2) 重要であるかないかという 1 つの観点に基づく分析では、文書同士の関わり合いを示すことができない。
- (3) どのような類の文書が必要であるかが分かっているにもかかわらず、すべての関連文書と一緒に表示されてしまう。

【0 0 1 0】

このようなシステムで、例えば、科学分野の文献を検索した場合、新しい技術の文献や、よく用いられるソフトウェアについての文献というような、主題として直接的な関連がないものがしばしば上位に表示されることがある。実際に参考にしたい文献はこのような文献ではない場合が多いが、従来のシステムでは、この問題を避けることができない。

【0 0 1 1】

本発明の課題は、ある文書がどのような理由で被参照文書を参照しているのかを同定し、同定された参照理由を利用して効率的に文書を検索する装置を提供することである。

【0 0 1 2】

【課題を解決するための手段】

図 1 は、本発明の装置の原理図である。本発明の第 1 の局面において、図 1 の装置は参照理由同定装置に対応し、被参照文書抽出手段 1 1、参照位置抽出手段 1 2、判断手段 1 3、および出力手段 1 4 を備える。

【0 0 1 3】

被参照文書抽出手段 1 1 は、与えられた文書データから被参照文書に関する文書情報を抽出する。参照位置抽出手段 1 2 は、文書データ内で被参照文書を参照している位置に関する情報を抽出する。判断手段 1 3 は、参照位置抽出手段 1 2 により抽出された情報を解析して、被参照文書が参照されている理由を判断する。そして、出力手段 1 4 は、被参照文書抽出手段 1 1 により抽出された情報と被

参照文書が参照されている理由とを含む出力情報を出力する。

【0014】

例えば、被参照文書は、与えられた文書データ内で参照されている参考文献等に対応し、被参照文書抽出手段11は、その文書データから被参照文書のタイトル、著者名等の文書情報を抽出する。また、参照位置抽出手段12は、文書データ内で参考文献の番号等を用いて被参照文書を参照している位置（参照位置）を探し、その位置が属する章、その位置の周辺の文字列等の情報を抽出する。

【0015】

また、判断手段13は、抽出された参照位置に関する情報に基づいて、被参照文書がその位置で参照されている理由（参照理由）を判断する。そして、出力手段14は、被参照文書の文書情報と参照理由とを対応付けて出力する。

【0016】

このような参照理由同定装置によれば、文書データから被参照文書の文書情報とその参照理由とが自動的に抽出され、利用者に提示される。したがって、利用者は、各被参照文書がどのような理由で参照元の文書により参照されているのかを認識することができ、参照理由から被参照文書の重要性を推定することもできる。さらに、利用者は、参照理由を利用して文書検索を行うこともできる。

【0017】

また、本発明の第2の局面において、図1の装置は情報検索装置に対応し、文書データベース手段15、参照関係格納手段16、検索手段17、および出力手段14を備える。

【0018】

文書データベース手段15は、文書データを格納する。参照関係格納手段16は、ある文書が被参照文書を参照している理由を表す参照理由を含む参照関係を格納する。検索手段17は、参照関係格納手段16に格納された参照関係を利用して、文書データベース手段15に格納された文書データの検索を行う。そして、出力手段14は、参照理由を含む検索結果を出力する。

【0019】

文書データベース手段15は、例えば、全文データベースに対応し、検索対象

となる文書データを格納する。参照関係格納手段 16 は、単にある文書が他の文書を参照しているという参照関係ではなく、その参照理由をも含む参照関係の情報を格納する。検索手段 17 は、格納された参照関係に含まれる参照理由を利用して、文書データベース手段 15 の文書データを検索する。そして、出力手段 14 は、検索に利用された参照理由に基づいて、得られた文書データの情報を出力する。

## 【0020】

このような情報検索装置によれば、検索された被参照文書が参照理由とともに利用者に提示されるため、利用者は、参照理由から被参照文書の重要性を推定することができる。さらに、利用者は、希望条件に合った適当な参照理由を指定して文書検索を行うこともでき、検索が効率化される。

## 【0021】

また、本発明の第3の局面において、図1の装置は文書分類装置に対応し、被参照文書抽出手段 11、参照位置抽出手段 12、判断手段 13、類似度判定手段 18、および出力手段 14 を備える。

## 【0022】

被参照文書抽出手段 11 は、与えられた文書データから被参照文書に関する文書情報を抽出する。参照位置抽出手段 12 は、文書データ内で被参照文書を参照している位置に関する情報を抽出する。判断手段 13 は、参照位置抽出手段 12 により抽出された情報を解析して、被参照文書が参照されている理由を判断する。類似度判定手段 18 は、複数の文書データ間において、被参照文書が参照されている理由を含む参照関係の類似度を計算し、それらの文書データを分類する。そして、出力手段 14 は、分類結果を出力する。

## 【0023】

類似度判定手段 18 は、例えば、2つの文書データ内で参照されている被参照文書の参照理由が与えられたとき、それらの参照理由を比較することで参照関係の類似度を計算する。このようにして、複数の文書データ内のすべての文書データ対について類似度が計算されると、それらの類似度に基づいて文書データを分類する。そして、出力手段 14 は、文書クラスタ等の情報を分類結果として出力

する。

【0024】

このような文書分類装置によれば、複数の文書データから被参照文書の文書情報とその参照理由とが自動的に抽出され、参照関係に基づいて複数の文書データが自動的に分類される。したがって、利用者は、参照理由を含む参照関係の観点から分類された文書群の情報を得ることができ、その情報を利用して文書検索を行うこともできる。

【0025】

例えば、図1の被参照文書抽出手段11は、後述する図2の文書構造解析部22と参考文献解析部25に対応し、図1の参照位置抽出手段12は、図2の文書構造解析部22と本文構造解析部24に対応し、図1の判断手段13は、図2の参照文脈解析部26に対応し、図1の出力手段14は、図2の対応判断部27と出力成形部28と出力インタフェース29に対応する。

【0026】

また、例えば、図1の文書データベース手段15は、後述する図11の全文データベース2に対応し、図1の参照関係格納手段16は、図11の参照関係データベース95に対応し、図1の検索手段17は、図11の検索装置92に対応する。

【0027】

また、例えば、図1の類似度判定手段18は、後述する図23の類似度判定装置152に対応する。

【0028】

【発明の実施の形態】

以下、図面を参照しながら、本発明の実施の形態を詳細に説明する。

一般に、システムが参照関係の特徴データを参照理由データ（関連性データ）としてあらかじめ持っていない場合、参照の有無を利用者に提示することは可能であるが、複数の参照理由を提示することは不可能か、あるいは非常に困難である。これに対して、システムが複数の参照理由データをあらかじめ持っている場合、参照の有無を利用者に提示することは簡単であり、複数の参照理由を提示す

ることも比較的簡単である。したがって、あらかじめ参照理由データをシステムに保持しておけば、それを利用して必要な情報を効率良く発見することができる。

#### 【 0 0 2 9 】

まず、参照理由の同定に関しては、参照している文書内における被参照文書の出現位置、引用方法、文章等を解析することで、被参照文書が参照されている理由を判断する。この解析においては、例えば、以下のような特徴付けが行われる。

( 1 ) 被参照文献は、参照している文献の知識を補うものである。例えば、被参照文献が条例であり、それを参照している文献が条例の追加である場合が、これに相当する。また、被参照文献が星の発見に関する文献であり、それを参照している文献がその星に関する新たなデータや仮説等に関する文献である場合も、これに相当する。

( 2 ) 被参照文献は、参照している文献の扱う分野をまとめたものである。このような被参照文献は、いわゆるレビュー論文のようなものに相当し、その分野の近年の進歩状況を照会するためによく引用される)

( 3 ) 参照している文献が被参照文献に対する反論を行っている。

( 4 ) 被参照文献は、参照している文献で紹介している人物の代表的な文献である。

#### 【 0 0 3 0 】

このようにして、数多くの文書进行处理することにより、参照されている文書がその分野でどのような位置付けにあるのかを知ることができる。

図 2 は、参照理由同定装置の構成図である。図 2 の参照理由同定装置は、入力インタフェース 2 1、文書構造解析部 2 2、書誌情報解析部 2 3、本文構造解析部 2 4、参考文献情報解析部 2 5、参照文脈解析部 2 6、対応判断部 2 7、出力成形部 2 8、および出力インタフェース 2 9 を備える。

#### 【 0 0 3 1 】

まず、入力インタフェース 2 1 は、文書データ 3 0 をテキストデータとして入力する。文書構造解析部 2 2 は、入力されたテキストデータを書誌情報、本文、



および参考文献の 3 つの部分に切り分けて、それぞれ書誌情報解析部 2 3、本文構造解析部 2 4、および参考文献解析部 2 5に渡す。

【0 0 3 2】

次に、書誌情報解析部 2 3 は、書誌情報からタイトル、著者名等の文献情報を抽出して出力し、本文構造解析部 2 4 は、本文から章構造と参考文献を引用している部分を抽出して出力し、参考文献解析部 2 5 は、参考文献の記載から人名、発行年、タイトル、誌名等の文献情報を抽出して出力する。

【0 0 3 3】

また、参照文脈解析部 2 6 は、本文中の参考文献を引用している部分を解析して参照理由を判断して出力し、対応判断部 2 7 は、本文中の参考文献を引用している部分と参考文献の記載から抽出された文献情報とを対応付けて出力する。

【0 0 3 4】

次に、出力成形部 2 8 は、書誌情報解析部 2 3 から出力された文献情報と、参照文脈解析部 2 6 から出力された参照理由と、対応判断部 2 7 から出力された対応関係と、参考文献解析部 2 5 から出力された参考文献の文献情報とをまとめて出力データ 3 1 を成形し、出力インタフェース 2 9 に渡す。

【0 0 3 5】

そして、出力インタフェース 2 9 は、出力データ 3 1 をディスプレイ画面等に出力する。出力データ 3 1 には、例えば、文書データ 3 0 の文献名と、参考文献および参照関係の組み合わせのリストが含まれる。

【0 0 3 6】

次に、図 3 から図 9 までを参照しながら、図 2 の参照理由同定装置の処理についてより詳細に説明する。

図 3 は、文書構造解析部 2 2 の処理のフローチャートである。パターンデータリスト 4 1、4 2、4 3、4 4 は、あらかじめ決められた文字列パターンの情報を含んでおり、不図示の記憶装置上に保持される。

【0 0 3 7】

文書構造解析部 2 2 は、まず、文書のテキストデータを読み込み（ステップ S 1）、パターンデータリスト 4 1 に格納された導入部パターン（はじめにパター

ン)を参照しながら、マッチするパターンを含む行を探し(ステップS2)、そのような行がテキストデータにあるか否かをチェックする(ステップS3)。文書が英語または日本語で作成されている場合、導入部パターンとしては、例えば、次のような文字列がパターンデータリスト41に格納される。

【0038】

1 Introduction

1. Introduction

Introduction

1 はじめに

1. はじめに

はじめに

1 背景

概要

このようなパターンにマッチする行があれば、マッチする最初の行の位置をP1として記憶する(ステップS4)。そして、パターンデータリスト42に格納されたキーワードパターンを参照しながら、その位置P1より前にキーワードパターンがあるか否かをチェックする(ステップS5)。キーワードパターンとしては、例えば、“\*\*\*”を任意のキーワードリストとして、次のような文字列がパターンデータリスト42に格納される。

【0039】

keyword \*\*\*

keyword: \*\*\*

キーワード \*\*\*

キーワード: \*\*\*

このようなパターンにマッチする行があれば、その位置を記憶する(ステップS6)。そして、パターンデータリスト43に格納されたアブストラクトパターンを参照しながら、記憶された位置より前にアブストラクトパターンがあるか否かをチェックする(ステップS7)。アブストラクトパターンとしては、例えば、次のようなパターン情報がパターンデータリスト43に格納される。

【0040】

Abstract

20語以上の文章で、文書の先頭の文章ではないもの。

このようなパターンにマッチする行があれば、その位置を記憶する（ステップS8）。そして、パターンデータリスト44に格納された参考文献パターンを参照しながら、マッチするパターンを含む行を探し（ステップS9）、そのような行が位置P1より後の部分にあるか否かをチェックする（ステップS10）。参考文献パターンとしては、例えば、次のような文字列がパターンデータリスト44に格納される。

【0041】

Reference

参考文献

参考文献

このようなパターンにマッチする行があれば、マッチする最後の行の位置をP2として記憶し（ステップS11）、記憶した位置P1、P2等を出力して（ステップS12）、処理を終了する。ステップS5においてキーワードパターンがない場合は、そのままステップS7以降の処理を行い、ステップS7においてアブストラクトパターンがない場合は、そのままステップS9以降の処理を行う。

【0042】

また、ステップS3において導入部パターンがない場合は、文書の先頭の位置をP1として記憶し（ステップS13）、ステップS9以降の処理を行う。また、ステップS10において参考文献パターンがない場合は、文書の末尾の位置をP2として記憶し（ステップS14）、ステップS12の処理を行う。

【0043】

ステップS12において出力される位置P1は書誌情報の記載と本文の記載の境界位置（切れ目）に対応し、位置P2は本文の記載と参考文献の記載の境界位置に対応する。

【0044】

次に、図4は、書誌情報解析部23の処理のフローチャートである。ストップ

ワードリスト 5 1 とパターンデタリスト 5 2 は、あらかじめ決められた文字列パターンの情報を含んでおり、記憶装置上に保持される。

【 0 0 4 5 】

書誌情報解析部 2 3 は、まず、テキストデータの先頭から P 1 までの部分を書誌情報として読み込み（ステップ S 2 1）、ストップワードリスト 5 1 を参照しながらテキストデータをチェックする。そして、最初にストップワードを含む文があればそれを読み飛ばし（ステップ S 2 2）、次の文をタイトルとして抽出する（ステップ S 2 3）。ストップワードとしては、タイトルではない文の先頭に現れる可能性のある文字列が用いられ、例えば、次のような文字列がストップワードリスト 5 1 に格納される。

【 0 0 4 6 】

解説

技術メモ

Technical Note

次に、パターンデタリスト 5 2 に格納された著者名パターンを参照しながら、タイトルとして記憶された文の次の文章からマッチするパターンを探し、それを抽出する（ステップ S 2 4）。著者名パターンとしては、例えば、図 4 に示すようなパターン情報がパターンデタリスト 5 2 に格納される。図 4 において、“〇〇”は辞書に登録された人名を表し、“————”は辞書に登録されていない未知語を表す。これらの著者名パターンは、次のように表すこともできる。

【 0 0 4 7 】

人名 + 未知語

未知語 + 未知語

アルファベット．人名

アルファベット．未知語

アルファベット．アルファベット．人名

アルファベット．アルファベット．未知語

次に、抽出されたデータをタイトルおよび著者名として記憶装置に書き出し（

ステップ S 25)、処理を終了する。

【0048】

次に、図 5 は、本文構造解析部 24 の処理のフローチャートである。パターンデータリスト 61、62 は、あらかじめ決められた文字列パターンの情報を含んでおり、記憶装置上に保持される。

【0049】

本文構造解析部 24 は、文書構造解析部 22 から渡された本文部分について、先頭から順に章構造を表すパターンを探していき、そのようなパターンの位置を記憶する処理を、それがなくなるまで繰り返す。次に、本文の先頭に戻って、参考文献の参照を表すパターンを探していき、そのようなパターンの位置を記憶する処理を、それがなくなるまで繰り返す。最後に、得られた章構造および参考文献の参照位置を出力する。

【0050】

本文構造解析部 24 は、まず、テキストデータの P1 から P2 までの部分を本文として読み込み（ステップ S 31）、パターンデータリスト 61 を参照しながら章構造を表すパターンを探していく（ステップ S 32）。章構造のパターンとしては、例えば、次のようなパターン情報がパターンデータリスト 61 に格納され、これらのパターンにマッチする行が探索される。

【0051】

数字 文字列（改行）

数字．文字列（改行）

ここで、パターンの先頭の数字は章番号を表しており、通常、本文の後にいくに従って大きくなる。また、文字の大きさや太さ等があらかじめ分かっているならば、それらの情報も参考にして章構造のパターンが探索される。

【0052】

次に、見つかったパターンに含まれる数字の増え方が適切か否かをチェックする（ステップ S 33）。そして、例えば、その数字が前の数字に繋がっていなければ、数字の増え方が適切でないと判断し、エラー処理を行って（ステップ S 34）、処理を終了する。

【 0 0 5 3 】

数字の増え方が適切であれば、見つかったパターンの位置を新たな章の先頭位置として記憶し（ステップ S 3 5）、まだ文章が残っているか否かをチェックする（ステップ S 3 6）。そして、文章が残っていれば、それがなくなるまでステップ S 3 2 以降の処理を繰り返す。

【 0 0 5 4 】

文章がなくなると、次に、位置 P 1 にポインタをセットし（ステップ S 3 7）、パターンデータリスト 6 2 を参照しながら、参考文献を参照していることを表すパターンを探す（ステップ S 3 8）。参考文献の参照を表すパターンとしては、例えば、次のようなパターン情報がパターンデータリスト 6 2 に格納される。

【 0 0 5 5 】

数字)

〔数字〕

〔文字列 年号〕

文字列 年号

これらのパターンにマッチする部分が見つかれば、その位置を参照位置として記憶し、文章がまだ残っているか否かをチェックする（ステップ S 3 9）。文章が残っていれば、それがなくなるまでステップ S 3 8 の処理を繰り返す。そして、文章がなくなると、得られた各章の先頭位置と参照位置を出力して（ステップ S 4 0）、処理を終了する。

【 0 0 5 6 】

次に、図 6 は、参考文献解析部 2 5 の処理のフローチャートである。パターンデータリスト 7 1 は、あらかじめ決められた文字列パターンの情報を含んでおり、記憶装置上に保持される。

【 0 0 5 7 】

参考文献解析部 2 5 は、文書構造解析部 2 2 から渡された参考文献部分を 1 文ずつ読み込み、参考文献パターンとマッチングして、マッチした人名、発行年、タイトル、誌名等の情報を順に記憶していく。このような処理を参考文献部分の行がなくなるまで繰り返す。

【0058】

参考文献解析部 25 は、まず、テキストデータの P2 から末尾までの部分を参考文献部分として読み込み（ステップ S41）、文章が残っているか否かをチェックする（ステップ S42）。文章が残っていれば、1 文を読み込み（ステップ S43）、その文とパターンデータリスト 71 の参考文献パターンとのマッチングを行い（ステップ S44）、マッチするパターンがあるか否かをチェックする（ステップ S45）。参考文献パターンとしては、例えば、次のようなパターン情報がパターンデータリスト 71 に格納される。

【0059】

人名 and 人名 年号 “タイトル” 誌名

〔人名 年号〕 人名 年号 “タイトル” 誌名

〔参照番号〕 人名 年号 “タイトル” 誌名

読み込んだ文がこのようなパターンにマッチすれば、その文に含まれる人名、年号、タイトル、誌名等の情報を参考文献情報として記憶し（ステップ S46）、ステップ S42 以降の処理を繰り返す。また、マッチするパターンがなければ、その文に含まれる情報を記憶せずに、ステップ S42 以降の処理を繰り返す。そして、ステップ S42 において文章がなくなると、得られた参考文献情報を出力して（ステップ S47）、処理を終了する。

【0060】

次に、図 7 は、参照文脈解析部 26 の処理のフローチャートである。参照文脈解析部 26 は、パターンデータリスト 71 は、あらかじめ決められた文字列パターンの情報を含んでおり、記憶装置上に保持される。

【0061】

参照文脈解析部 26 は、言語解析部 81 と参照特徴－参照理由対応表 82 を含み、本文構造解析部 24 から渡された章構造および参照位置の情報を解析して、参照理由を判断する。この参照特徴－参照理由対応表 82 は、あらかじめ人手または学習により作成される。参照文脈解析部 26 は、各参照位置について、その位置を含む章の情報と、その位置の周辺の文字列を言語解析部 81 により解析した結果と、参考文献の引用パターンの 3 つの情報に基づき、参照特徴－参照理由

対応表 82 を用いて参照理由を判断する。

【0062】

参照文脈解析部 26 は、まず、本文と章構造を読み込み（ステップ S51）、すべての参照位置について処理が完了したか否かをチェックする（ステップ S52）。処理が完了していなければ、次の参照位置を読み込み（ステップ S53）、章構造を参照しながらその位置が属する章の番号とその位置の周辺の文字列を抽出する（ステップ S54）。

【0063】

次に、言語解析部 81 は、抽出された文字列の形態素解析、構文解析、意味解析等を行って、参照位置周辺の特徴を抽出する（ステップ S55）。例えば、形態素解析については、以下の文献〔1〕の pp. 117-137 のアルゴリズムが用いられ、構文解析については、この文献の pp. 140-199 のアルゴリズムが用いられる。また、意味解析については、この文献の pp. 200-231 のアルゴリズムが用いられる。

〔1〕長尾真（ながおまこと）、「自然言語処理」，岩波書店，1996

次に、参照文脈解析部 26 は、参照特徴-参照理由対応表 82 を参照しながら、参照理由を判断する（ステップ S56）。参照理由としては、例えば、以下のようなものが考えられる。

（1）反論（answer）

参考文献の内容と異なる意見を記述するために、参考文献を引用する。

（2）応用（application）

応用分野を紹介するために、参考文献を引用する。

（3）基礎（basic）

基礎的な研究や先人の仕事を紹介するために、参考文献を引用する。

（4）対立する意見（contraposition）

特定の内容と対立する意見や対照的な意見を紹介するために、参考文献を引用する。



(5) 人物の代表的な文献 (human)

特定の人物の代表的な文献を紹介するために、参考文献を引用する。

(6) 関連のある仕事 (related work)

特定の内容と関連のある仕事を紹介するために、参考文献を引用する。

(7) まとめ (review)

特定の分野のまとめを紹介するために、参考文献を引用する。

(8) 使用ソフトウェア (software)

シミュレーション等に使用したソフトウェアを紹介するために、参考文献を引用する。

(9) 使用技術 (technique)

実験、シミュレーション等に使用した技術を紹介するために、参考文献を引用する。

(10) 弱い関係 (weak correlation)

特定の内容と弱い関係を持つような内容を紹介するために、参考文献を引用する。

(11) 類似 (similar)

特定の内容と類似した内容を紹介するために、参考文献を引用する。

【0064】

参照特徴－参照理由対応表 8 2 には、このような参照理由のカテゴリと参照特徴との対応関係が格納されている。この参照特徴は、対応する参照理由を表現する本文の記述方法の特徴を表し、参考文献の参照位置を含む章の番号、参照位置の周辺の文字列の文脈、参考文献の引用パターン等の情報を含む。上述の参照理由の場合は、“\*\*”を参考文献を示す文字列として、例えば、以下のような対応関係が参照特徴－参照理由対応表 8 2 に格納される。

【0065】

参照位置 = 1 章、かつ、言語解析結果 = 否定表現 → 反論

文脈 = “. . . \*\* というのが通説であるが、しかしながら、本研究では. .

.” → 反論

文脈 = “This approach is used in \*\*...” → 応用

参照位置 = 1 章または 2 章 → 基礎

文脈 = “The first idea...due to \*\*” または “In previous research...\*\*”  
→ 基礎

文脈 = “Unlike previous...\*\*, new...” → 対立する意見

参照位置 = 1 章、かつ、文脈 = “. . . が提案された \*\*” または “. . . が提案されている \*\*” → 関連のある仕事

文脈 = “...is reviewd in \*\* ” または “\*\* reviewd” または “see \*\* for a n overview” → まとめ

文脈 = “We use...similar to \*\*” または “see \*\* for a similar...approach” または “...is also implemented in \*\*” → 類似

参照文脈解析部 2 6 は、ステップ S 5 4 で抽出された情報やステップ S 5 5 で得られた解析結果を参照特徴として参照特徴－参照理由対応表 8 2 を検索し、対応する参照理由を取得し、それを記憶する（ステップ S 5 7）。そして、ステップ S 5 2 以降の処理を繰り返し、ステップ S 5 2 においてすべての参照位置についての処理が完了すると、各参照位置毎に参照理由を出力して（ステップ S 5 8）、処理を終了する。

【 0 0 6 6 】

このように、参照理由は、参照を行っている文書中での被参照文書に関する表現の特徴に基づいて決められる。この参照理由を判断する方法としては、上述のような参照特徴－参照理由対応表 8 2 を保持する方法のほかに、統計的な方法を用いることもできる。このような方法は機械学習とも呼ばれる。この場合、参照特徴としては、次のようなものが用いられる。

- (1) 参考文献の出現回数
- (2) 参考文献の参照位置：文書の先頭からの距離、出現した章
- (3) 参照位置の周辺情報：共起している単語、よく現れるフレーズ

これらの参照特徴に対して、被参照文書がどのような理由で参照されているのかを人間があらかじめ正解として用意しておき、機械もしくは人手で抽出した参照特徴を参照理由に対応付ける。このとき、参照特徴と参照理由のセットを複数用意しておき、統計的な方法により参照特徴から参照理由への対応を得る。この

統計的な方法としては、例えば、次の文献〔2〕の p p. 5 2 5 - 6 5 2 に記載された決定木、ニューラルネット、最近傍法、ベイズ推定等の方法が挙げられる。

〔2〕 S. Russell and P. Norvig, 「エージェントアプローチ 人工知能」, 共立出版, 1 9 9 7

次に、具体的な文書を例に挙げて、参照理由同定装置の動作を説明する。図 8 は、“Fujino”の文献を参考文献として参照している“Jpn.J.Appl.Phys.”の文書を示している。

【0 0 6 7】

まず、この文書が入力インタフェース 2 1 を通して入力されると、入力されたデータは、文書構造解析部 2 2 により、書誌情報（●Jpn.J.Appl.Phys.～KEYWORD;...）と本文（1 Introduction～）と参考文献（Reference～K.Fujino,...）の 3 つの部分に分割される。

【0 0 6 8】

ここでは、図 3 の分割位置 P 1 は 1 章の“1 Introduction”の行の直前に設定され、分割位置 P 2 は“Reference”の行の直前に設定されている。したがって、“1 Introduction”の前の行までが書誌情報に対応し、1 章から 3 章までが本文に対応し、“Reference”の行以降が参考文献の部分に対応する。

【0 0 6 9】

このうち、書誌情報は書誌情報解析部 2 3 により解析され、この文書の著者が Koji Tsukamoto らであること、タイトルが“Morphology Evolution....”であること等が同定される。また、本文は本文構造解析部 2 4 により解析され、章構造が認識されるとともに、本文中で参考文献を引用している部分として、“6,8,10-13)”、“8,11,13)”のようなパターンが認識される。

【0 0 7 0】

このような認識結果は参照文脈解析部 2 6 に送られ、参考文献がどのような章のどのような文章の中で参照されているかに基づいて、参照理由が判断される。

例えば、文献6)が1章で4回参照され、3章で2回参照されているものとする  
と、この文献は相対的に非常に多く参照されていることになる。このことから、  
文献6)は重要な文献であり、図8の文書の基礎となる文献であることが分かる  
。

## 【0071】

また、文献6), 8), 10-13)は“Some methods have been.....However”  
という表現で参照されていることから、この文書とは対立する意見を述べた  
文献であることが分かる。また、文献8), 11), 13)は“Similar results  
were also reported.....”という表現で参照されていることから、この文書  
と類似した文献であることが分かる。

## 【0072】

また、参考文献の部分は参考文献解析部25により解析され、参考文献の文献  
情報が抽出される。例えば、文献6)の作者は“Fujino”であり、発行年は“1  
991年”であり、掲載誌は“J.Electrochem.Soc.”であること等が認定される  
。

## 【0073】

次に、対応判断部27により、文献6)の参照を表している表記“6)”と文献  
6)の文献情報とが対応付けられる。そして、出力成形部28により、書誌情報  
解析部23、参照文脈解析部26、対応判断部27、および参考文献解析部25  
から出力された情報がまとめられて、出力インタフェース29から出力される。

## 【0074】

このようにして同定された参照理由を用いれば、ある文書が他の文献を参照し  
ている理由を利用して文書を検索することにより、検索を効率化することができ  
る。

## 【0075】

例えば、参照理由を用いない場合の参照関係のデータ構造は、図9に示すよう  
になる。ここで、各アルファベット文字は1つの文書を表し、矢印の元の文書は  
矢印の先の文書を参照している。これに対して、参照理由を用いた場合の参照関  
係のデータ構造は、図10に示すようになる。図10では、文書間の参照／被参

照の情報だけでなく、参照理由毎の参照／被参照の情報が示されている。

【0076】

図27に示した従来の情報検索システムにおいて、参照関係データベースに図10のような参照理由を用いた参照／被参照関係を格納した場合、システム構成は図11のようになる。

【0077】

図11の情報検索システムは、入力装置91、検索装置92、表示装置93、全文データベース94、および参照関係データベース95を備え、検索装置92は、制御部101、検索部102、および選択部103を含む。参照関係データベース95には、図2の参照理由同定装置により抽出された情報が格納される。

【0078】

入力装置91によりある文書が指定されると、制御部101による制御に基づいて、検索部102は、全文データベース94の文書を検索する。そして、選択部103は、参照関係データベース95の参照関係を用いて、指定された文書と共引用関係にある文献や似たような参照を行っている文献を、関連文書として選択し、表示装置93は、それらの関連文書を表示する。

【0079】

図12は、このような検索結果の表示例を示している。参照関係データベース95に指定された文書の参照関係111が格納されており、文書X、W、V等が関連文書として選択されたとする。このとき、表示装置93の画面には検索結果112が表示され、各文書に該当する参照理由の欄にマーク“○”が記される。ここでは、関連文書が参照されている理由として、“answer”、“application”、“basic”等が提示されている。

【0080】

図13は、図11のシステムによる検索結果の表示処理のフローチャートである。選択部103は、まず、参照関係データベース95から指定された文書の参照関係データを読み込み（ステップS61）、そのデータに記述された参照理由を読み込む（ステップS62）。次に、指定された文書が参照している被参照文書の文献情報を読み込んで、図12のような表形式の検索結果を表示し（ステッ

プ S 6 3)、該当する参照理由の欄にマークを付ける(ステップ S 6 4)。

【0081】

次に、参照関係データが残っているか否かをチェックし、データが残っていれば、ステップ S 6 1以降の処理を繰り返す。そして、指定された文書のすべての参照関係データが処理されると、表示処理を終了する。

【0082】

このような表示処理によれば、ある文書に関連する重要文書を表示する際、図 28に示したような単一の基準(参照回数)に基づく順序だけでなく、その重要文書が参照されている理由も明示される。したがって、参照理由に基づいて検索結果をさらに絞り込むことも可能になり、文書検索が効率化される。

【0083】

また、関連文書を単に重要度順に表示するだけでなく、図 14に示すように、関連文書を参照理由毎に分類して重要度順に表示することも可能である。この場合、各文書の重要度は、上記単一の基準に基づく重要度に加え、以下のような要素を加味して決定される。

- (1) それぞれの参照理由
- (2) 対応する参照理由により参照された回数
- (3) 参照元の文書の重要度
- (4) 出典の重要度
- (5) 利用者に参照された回数

表示された参照理由の中から利用者が適当なものを選択して、さらに検索を指示した場合、検索部 102は、参照関係データベース 95を参照しながら全文データベース 94を検索する。そして、指定された参照理由で参照されている文書が関連文書として表示される。

【0084】

このように、図 11の情報検索システムによれば、検索結果に参照理由を付与したり、参照理由を用いて検索を行ったりすることで、利用者が情報検索の効率を改善することができる。

【0085】

次に、被参照文書に関する情報を抽出して提示する情報提示装置について説明する。図15は、このような情報提示装置の構成図である。図15の情報提示装置は、図11の情報検索システムの構成要素に加えて、入力インタフェース121、検索入力インタフェース122、検索出力インタフェース123、参照理由同定装置124、および参照関係変換装置125をさらに備える。

#### 【0086】

入力インタフェース121は、ある文書を参照している複数の文書（文書群）を入力し、全文データベース94に格納するとともに、参照理由同定装置124に与える。参照理由同定装置124は、図2のような構成を持ち、与えられた文書の参照関係を求めて、参照関係変換装置125に出力する。

#### 【0087】

参照関係変換装置125は、与えられた参照関係を変換して、参照関係データベース95に格納する。例えば、図16に示すように、文書Aが文書B、C、Dをそれぞれb、c、dという参照理由で参照しているという参照関係が与えられたとする。このとき、参照関係変換装置125は、図17に示すように、文書B、C、Dが文書Aからそれぞれb、c、dという理由で参照されていることを表すデータを加えて、図18に示すような参照関係データを生成する。そして、このデータを参照関係データベース95に格納する。

#### 【0088】

利用者は、検索入力インタフェース122から、キーワードや文書を指定して検索を行う。そして、検索出力インタフェース123により、検索結果の文書の情報が表示されると、その文書の参照理由を用いてさらに検索を行う。

#### 【0089】

検索入力インタフェース122および検索出力インタフェース123では、以下のような入出力が行われる。

- (1) 利用者が文書を指定することで、その文書の本文を表示する。
- (2) 興味のある参照関係をフィルタとして使用する。
- (3) 参照理由を明示的に表示して参照関係を表示する。
- (4) GUI (graphical user interface) を用いて、参照理由を明示しながら

、あるいは参照関係を表す線または矢印の種類を参照理由によって変えながら、参照関係を画像表示する。

#### 【0090】

例えば、複数の文書を発行年順に並べて、文書間の参照関係を矢印で表すと、図19のようになる。ここでは、円形のマークが1つの文書を表し、複数の文書が時間軸に沿って時系列に表示されている。また、参照元の文書と被参照文書は参照理由毎に異なる線種の矢印で結ばれる。

#### 【0091】

このようなGUIを用いれば、複雑な参照関係における各文書の位置付けを明確に認識することができ、文書検索が効率化される。さらに、この表示画面中の文書マークを利用者が指定することにより、対応する文書の内容を表示したり、その文書が参照している文書のリストを表示したりすることもできる。

#### 【0092】

図19では、矢印の線種により参照理由を区別しているが、矢印の色によりこれを区別してもよく、参照理由を文字列で表示してもよい。

次に、参照関係を利用したキーワード抽出装置について説明する。図20は、このようなキーワード抽出装置の構成図である。図20のキーワード抽出装置は、図2の参照理由同定装置の構成要素に加えて、キーワード抽出部131をさらに備える。

#### 【0093】

本文構造解析部24は、解析結果をキーワード抽出部131に出力し、キーワード抽出部131は、本文中で参考文献が参照されている位置の周辺の文章からキーワードを抽出し、それを参考文献のキーワードとして出力する。この場合、出力インタフェース29から出力される出力データ132には、文書データ30の文献名に加えて、参考文献、参照関係、およびキーワードの組み合わせのリストが含まれる。

#### 【0094】

このようなキーワード抽出装置によれば、文献の著者自身が選択したキーワードだけでなく、その文献を参照している他の文書に基づくキーワードも付与され



ることになる。したがって、主観的なキーワードだけではなく、客観的なキーワードを付与することが可能になる。

【0095】

図21は、このような参照関係を利用したキーワードの例を示している。ここでは、ある文献の著者が“Machine Learning”、“Decision Tree”等をキーワードとして付与し、キーワード抽出装置は、その文献を参照している他の文書の文脈から“Corpus”という新たなキーワードを抽出してキーワードのデータ構造に付加している。

【0096】

図22は、このようにして付与されたキーワードを用いて文書を検索する情報検索システムを示している。図22の情報検索システムは、図11に示したシステムの構成要素に加えて、文献キーワード記憶部141をさらに備える。文献キーワード記憶部141には、文献にあらかじめ付与されたキーワードが格納され、参照関係データベース95には、図20のキーワード抽出装置により抽出された情報が格納される。

【0097】

検索部102は、文献キーワード記憶部141と参照関係データベース95のキーワードを参照しながら全文データベース94の文書を検索し、選択部103は、文献キーワード記憶部141と参照関係データベース95の情報をを用いて関連文書を選択する。このとき、表示の優先順位は、キーワードの一致度、参照関係に基づく重要度、および文書に対するアクセス回数を加味して決められる。

【0098】

このような情報検索システムによれば、あらかじめ付与されたキーワードに加えて、参照関係を用いて客観的に付与されたキーワードを利用して検索することができ、検索効率が向上する。

【0099】

また、文書間の参照関係を利用して多数の文書を分類することも可能である。図23は、このような参照関係に基づく分類を行う文書分類装置の構成図である。図23の文書分類装置は、参照関係データベース95、キーワード抽出装置1

51、および類似度判定装置152を備える。

【0100】

キーワード抽出装置151は、図20のような構成を持ち、複数の文書データ153から参考文献に関する参照関係、キーワード等の情報を抽出して、参照関係データベース95に格納する。類似度判定装置152は、入力された文書間で参照理由を含む参照関係を比較し、参照関係の類似度に基づくクラスタリングを行って、文書クラスタ154を出力する。

【0101】

例えば、文書aと文書bの間で参照関係の類似度を計算する際に、次のような計算式が用いられる。

【0102】

【数1】

$$\text{sim}(a, b) = \frac{1}{\sqrt{n_a} \sqrt{n_b}} \sum_i^{n_a} \sum_j^{n_b} \delta(r_{ai}, r_{bj})$$

【0103】

この計算式において、 $\text{sim}(a, b)$ は、文書aと文書bの間の類似度を表す。また、 $n_a$ 、 $n_b$ は、それぞれ、文書a、文書bが参照している被参照文書の数であり、 $r_{ai}$ 、 $r_{bj}$ は、それぞれ、文書a、文書bによる参照を表すベクトルであり、被参照文書とその参照理由の属性を持つ ( $i = 1, 2, \dots, n_a$ ,  $j = 1, 2, \dots, n_b$ )。

【0104】

$\delta(r_{ai}, r_{bj})$ は、 $r_{ai}$ と $r_{bj}$ の間の類似度を表す関数である。この関数は、例えば、 $r_{ai}$ と $r_{bj}$ が同じ被参照文書を同じ参照理由で参照している場合は $\delta(r_{ai}, r_{bj}) = 1$ となり、そうでない場合は $\delta(r_{ai}, r_{bj}) = 0$ となるように定義される。また、 $r_{ai}$ と $r_{bj}$ が同じ被参照文書を異なる参照理由で参照していた場合に、 $\delta(r_{ai}, r_{bj})$ に1または0.5を割り当てるという方法もある。

【0105】

被参照文書についても同じ方法で類似度計算が行われ、計算結果に基づいて文書のクラスタリングが行われる。文書のクラスタリングには、例えば、上述の文献〔1〕の pp. 436-438 に記載されているアルゴリズムが用いられる。これにより、参照関係が互いに類似している複数の文書が同じクラスに分類される。また、このような類似度計算に、さらにキーワードの一致度、単語の出現頻度等を加味したクラスタリングを行うこともできる。

## 【0106】

このように、文書間の参照関係に基づいて文書を分類することにより、利用者は、参照理由を含む参照関係の観点から分類された文書群の情報を得ることができる。また、その分類結果を利用して文書を検索すれば、検索が効率化される。

## 【0107】

図24は、図19に示した参照関係の時系列表示において、文書の分類結果を利用した例を示している。図24において、縦軸は、分類により得られた文書の種類を表し、横軸は、時間を表す。このような表示方法によれば、同じような参照関係を持つ文書が近接して配置され、多数の文書間の参照関係をより分かりやすく表示することができる。

## 【0108】

以上説明した実施形態においては、文書間の参照関係に関する処理を行っているが、同様にして、任意の情報間における参照関係を処理することもできる。例えば、参考文献の代わりに、他のテキストデータ、画像データ、音声データ、プログラムリスト等が参照されている場合、それらの情報の参照理由が同定され、情報検索に利用される。

## 【0109】

ところで、図2の参照理由同定装置、図11および図22の情報検索システム、図15の情報提示装置、図20のキーワード抽出装置、および図23の文書分類装置は、図25に示すような情報処理装置（コンピュータ）を用いて構成することができる。図25の情報処理装置は、CPU（中央処理装置）161、メモリ162、入力装置163、出力装置164、外部記憶装置165、媒体駆動装置166、およびネットワーク接続装置167を備え、それらはバス168によ

り互いに接続されている。

【0 1 1 0】

メモリ 1 6 2 は、例えば、ROM (read only memory)、RAM (random access memory) 等を含み、処理に用いられるプログラムとデータを格納する。CPU 1 6 1 は、メモリ 1 6 2 を利用してプログラムを実行することにより、必要な処理を行う。

【0 1 1 1】

入力装置 1 6 3 は、例えば、キーボード、ポインティングデバイス、タッチパネル等であり、利用者からの指示や情報の入力に用いられる。出力装置 1 6 4 は、例えば、ディスプレイ、プリンタ、スピーカ等であり、利用者へのメッセージや処理結果の出力に用いられる。

【0 1 1 2】

外部記憶装置 1 6 5 は、例えば、磁気ディスク装置、光ディスク装置、光磁気ディスク (magneto-optical disk) 装置等であり、図 1 1 の全文データベース 9 4、参照関係データベース 9 5、および図 2 2 の文献キーワード記憶部 1 4 1 として用いられる。また、情報処理装置は、この外部記憶装置 1 6 5 に、上述のプログラムとデータを保存しておき、必要に応じて、それらをメモリ 1 6 2 にロードして使用することができる。

【0 1 1 3】

媒体駆動装置 1 6 6 は、可搬記録媒体 1 6 9 を駆動し、その記録内容にアクセスする。可搬記録媒体 1 6 9 としては、メモリカード、フロッピーディスク、CD-ROM (compact disk read only memory)、光ディスク、光磁気ディスク等、任意のコンピュータ読み取り可能な記録媒体が用いられる。利用者は、この可搬記録媒体 1 6 9 に上述のプログラムとデータを格納しておき、必要に応じて、それらをメモリ 1 6 2 にロードして使用することができる。

【0 1 1 4】

ネットワーク接続装置 1 6 7 は、任意のネットワーク (回線) を介して外部の装置と通信し、通信に伴うデータ変換を行う。情報処理装置は、必要に応じて、ネットワーク接続装置 1 6 7 を介して上述のプログラムとデータを外部の装置か

ら受け取り、それらをメモリ 1 6 2 にロードして使用することができる。

【0 1 1 5】

図 2 6 は、図 2 5 の情報処理装置にプログラムとデータを供給することのできるコンピュータ読み取り可能な記録媒体を示している。可搬記録媒体 1 6 9 や外部のデータベース 1 7 0 に保存されたプログラムとデータは、メモリ 1 6 2 にロードされる。そして、CPU 1 6 1 は、そのデータを用いてそのプログラムを実行し、必要な処理を行う。

【0 1 1 6】

【発明の効果】

本発明によれば、文書等の情報内で他の情報を参照している部分を解析することにより、その情報が参照されている理由を判断することができる。また、情報の参照理由を利用者に提示すれば、利用者は複数の参照理由から希望するものを選択することができ、情報検索が効率化される。さらに、参照理由を用いて情報間の参照関係を提示することにより、各情報の位置付けが明確に認識されるようになる。

【図面の簡単な説明】

【図 1】

本発明の装置の原理図である。

【図 2】

参照理由同定装置の構成図である。

【図 3】

文書構造解析部の処理のフローチャートである。

【図 4】

書誌情報解析部の処理のフローチャートである。

【図 5】

本文構造解析部の処理のフローチャートである。

【図 6】

参考文献解析部の処理のフローチャートである。

【図 7】

参照文脈解析部の処理のフローチャートである。

【図 8】

文書の例を示す図である。

【図 9】

第 1 の参照関係を示す図である。

【図 1 0】

第 2 の参照関係を示す図である。

【図 1 1】

情報検索システムの構成図である。

【図 1 2】

第 1 の検索結果表示を示す図である。

【図 1 3】

表示処理のフローチャートである。

【図 1 4】

第 2 の検索結果表示を示す図である。

【図 1 5】

情報提示装置の構成図である。

【図 1 6】

第 3 の参照関係を示す図である。

【図 1 7】

第 4 の参照関係を示す図である。

【図 1 8】

第 5 の参照関係を示す図である。

【図 1 9】

第 1 の時系列表示を示す図である。

【図 2 0】

キーワード抽出装置の構成図である。

【図 2 1】

参照関係を用いたキーワードを示す図である。

【図 2 2】

キーワードを用いた検索を示す図である。

【図 2 3】

文書分類装置の構成図である。

【図 2 4】

第 2 の時系列表示を示す図である。

【図 2 5】

情報処理装置の構成図である。

【図 2 6】

記録媒体を示す図である。

【図 2 7】

従来の情報検索システムを示す図である。

【図 2 8】

従来の表示形式を示す図である。

【符号の説明】

- 1、9 2 検索装置
- 2、9 4 全文データベース
- 3、9 5 参照関係データベース
- 4、1 0 1 制御部
- 5、1 0 2 検索部
- 6、1 0 3 選択部
- 1 1 被参照文書抽出手段
- 1 2 参照位置抽出手段
- 1 3 判断手段
- 1 4 出力手段
- 1 5 文書データベース手段
- 1 6 参照関係格納手段
- 1 7 検索手段
- 1 8 類似度判定手段

- 2 1、1 2 1 入力インタフェース
- 2 2 文書構造解析部
- 2 3 書誌情報解析部
- 2 4 本文構造解析部
- 2 5 参考文献解析部
- 2 6 参照文脈解析部
- 2 7 対応判断部
- 2 8 出力成形部
- 2 9 出力インタフェース
- 3 0、1 5 3 文書データ
- 3 1、1 3 2 出力データ
- 4 1、4 2、4 3、4 4、5 2、6 1、6 2、7 1 パターンデータリスト
- 5 1 ストップワードリスト
- 8 1 言語解析部
- 8 2 参照特徴－参照理由対応表
- 9 1、1 6 3 入力装置
- 9 3 表示装置
- 1 1 1 参照関係
- 1 1 2 検索結果
- 1 2 2 検索入力インタフェース
- 1 2 3 検索出力インタフェース
- 1 2 4 参照理由同定装置
- 1 2 5 参照関係変換装置
- 1 3 1 キーワード抽出部
- 1 4 1 文献キーワード記憶部
- 1 5 1 キーワード抽出装置
- 1 5 2 類似度判定装置
- 1 5 4 文書クラスタ
- 1 6 1 C P U



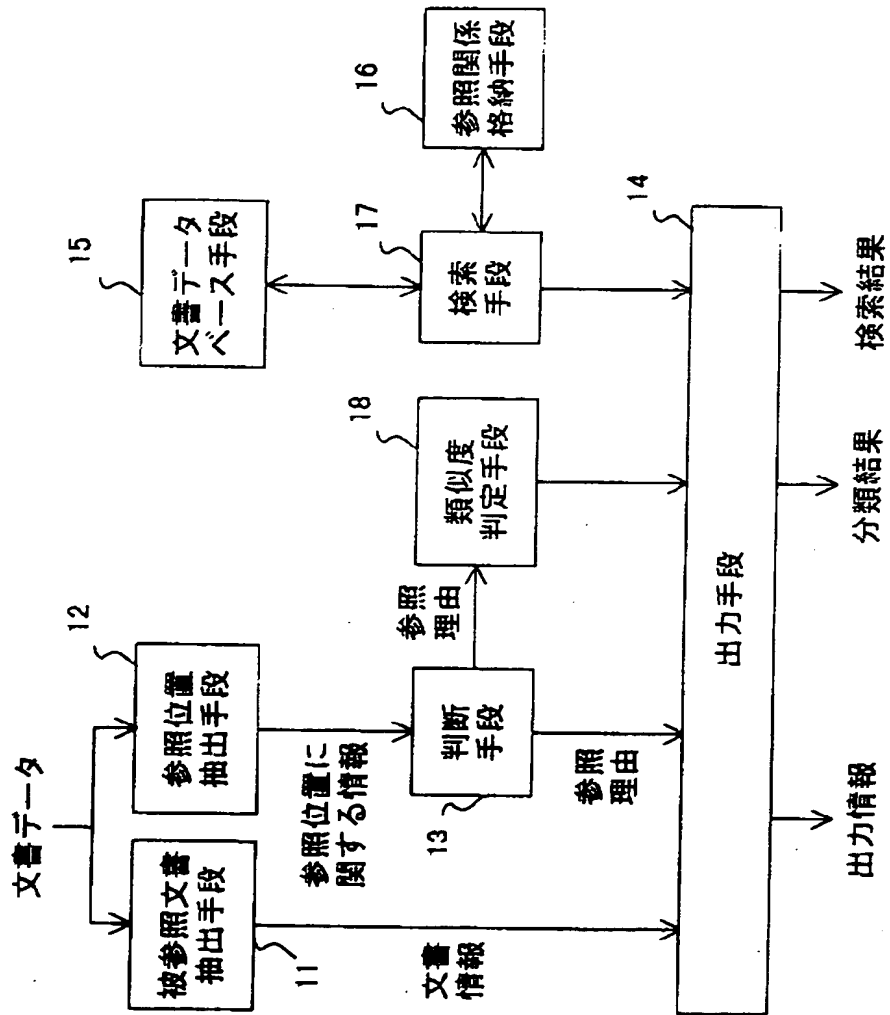
- 1 6 2    メモリ
- 1 6 4    出力装置
- 1 6 5    外部記憶装置
- 1 6 6    媒体駆動装置
- 1 6 7    ネットワーク接続装置
- 1 6 8    バス
- 1 6 9    可搬記録媒体
- 1 7 0    データベース

【書類名】

図面

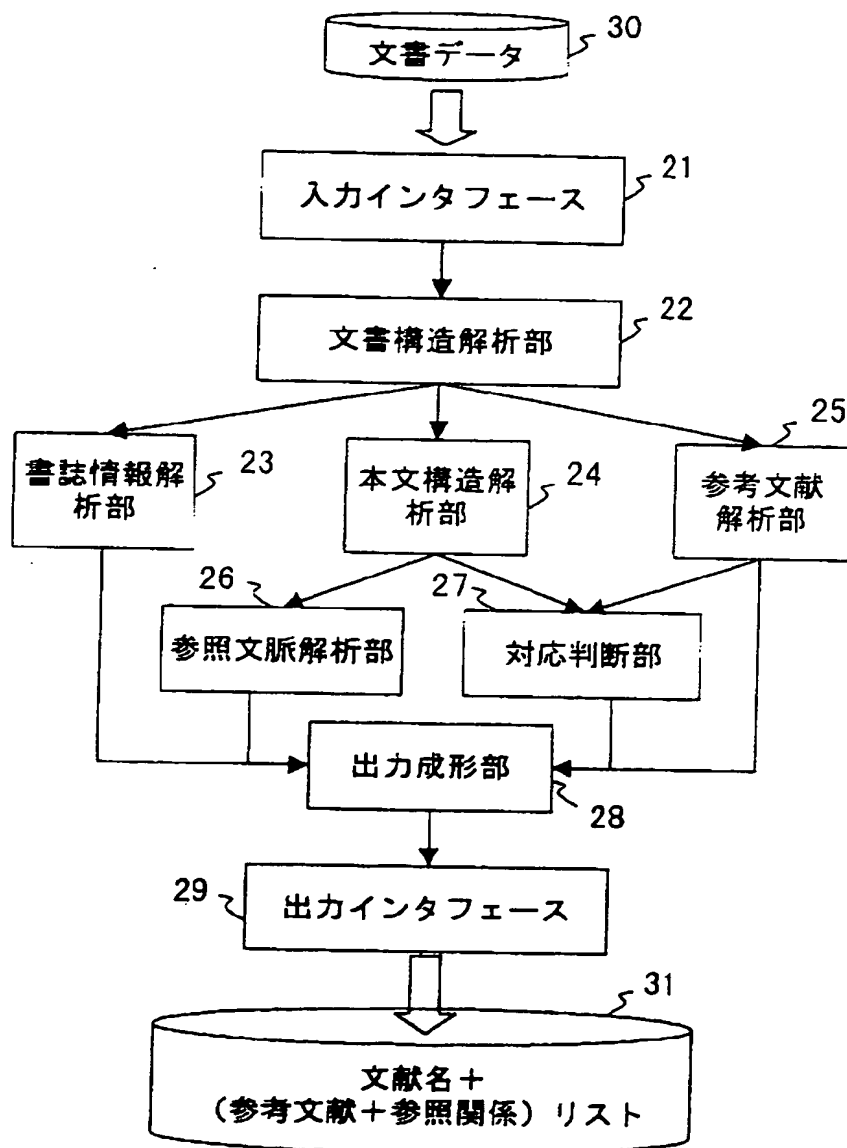
【図 1】

# 本 発 明 の 原 理 図



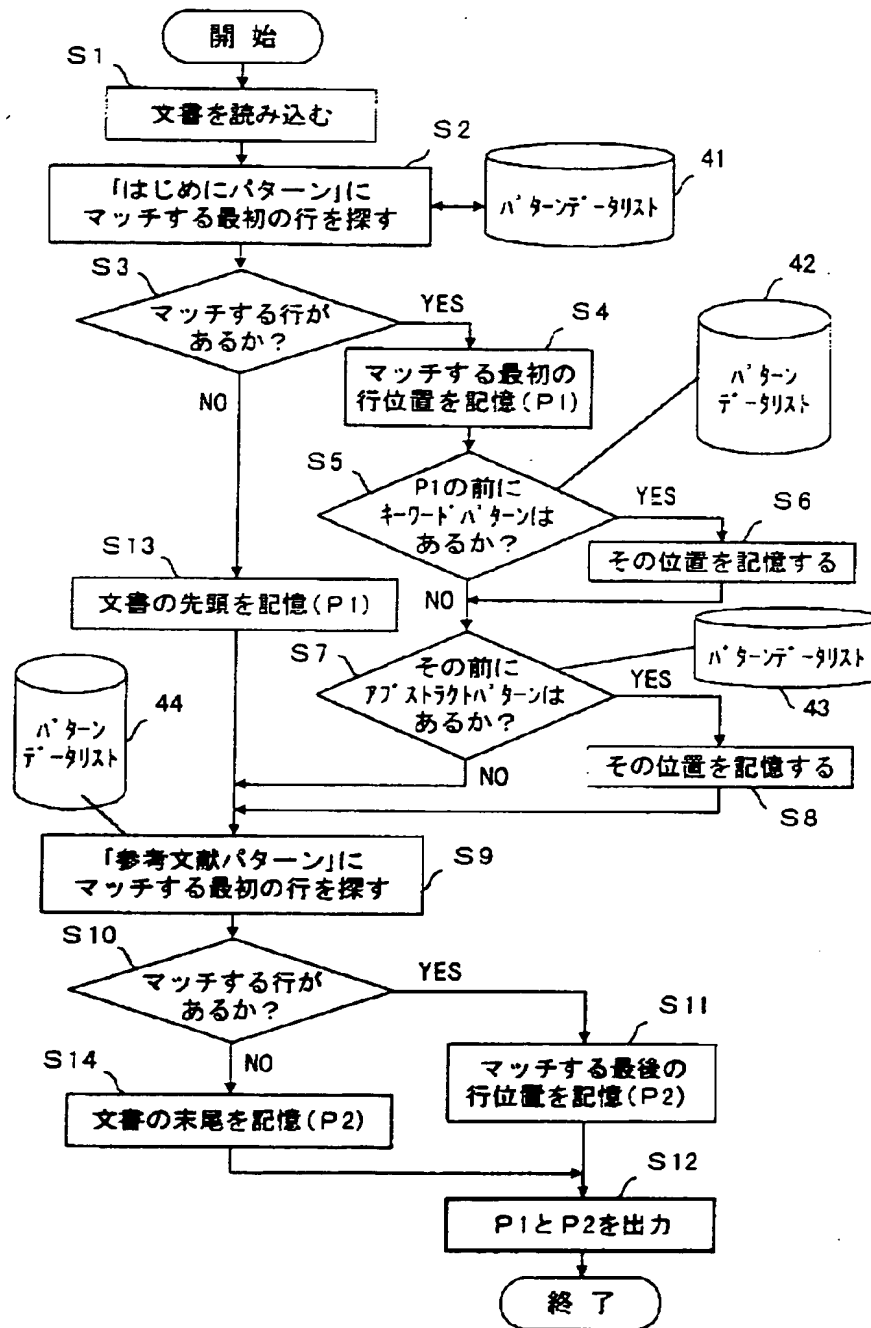
【図 2】

参 照 理 由 同 定 装 置 の 構 成 図



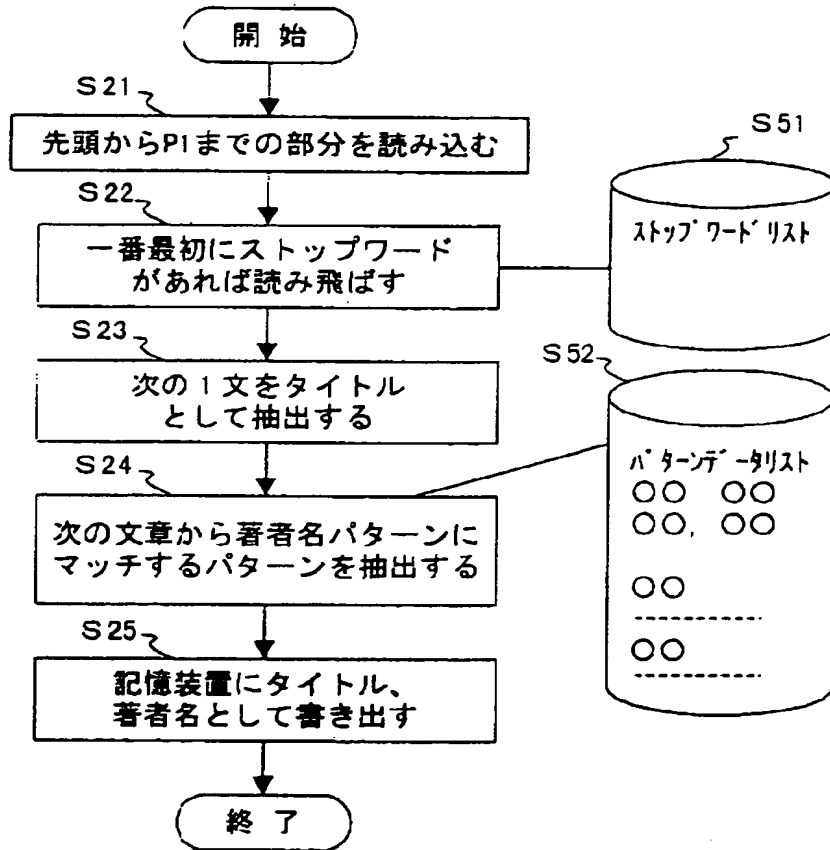
【図 3】

文書構造解析部の処理のフローチャート



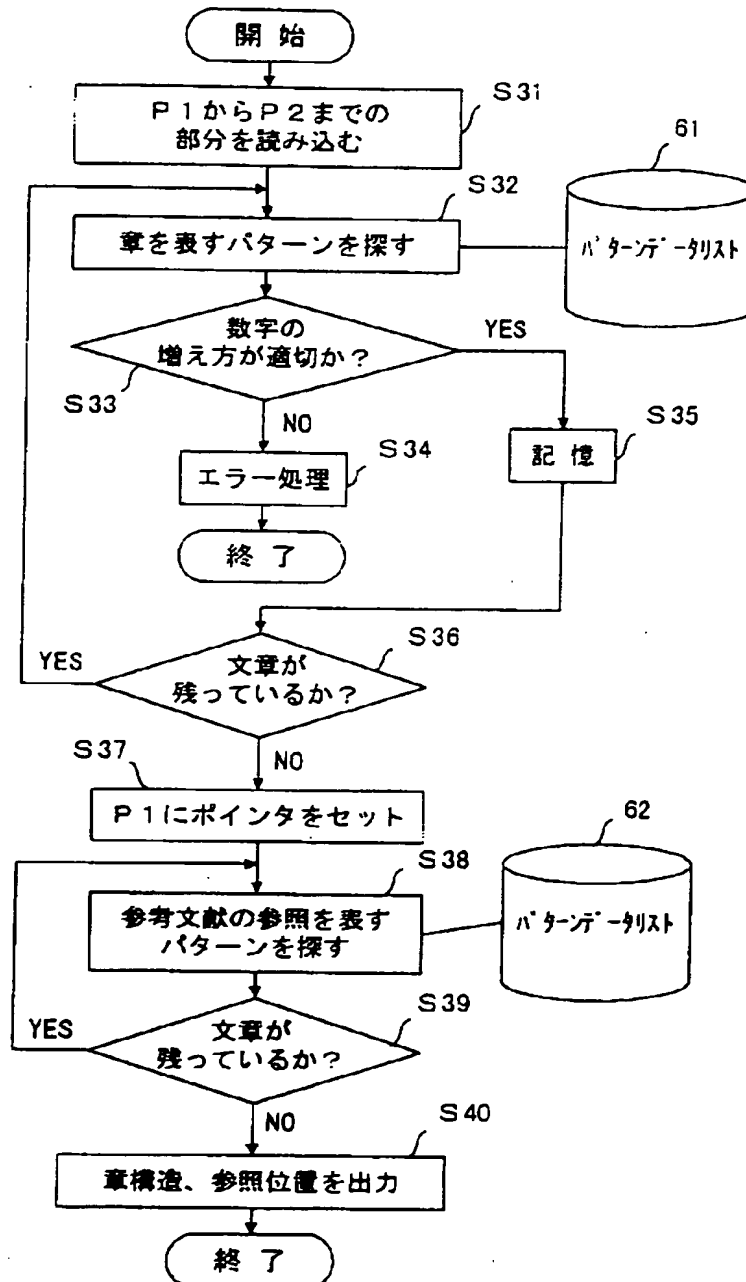
【図 4】

書誌情報解析部の処理のフローチャート



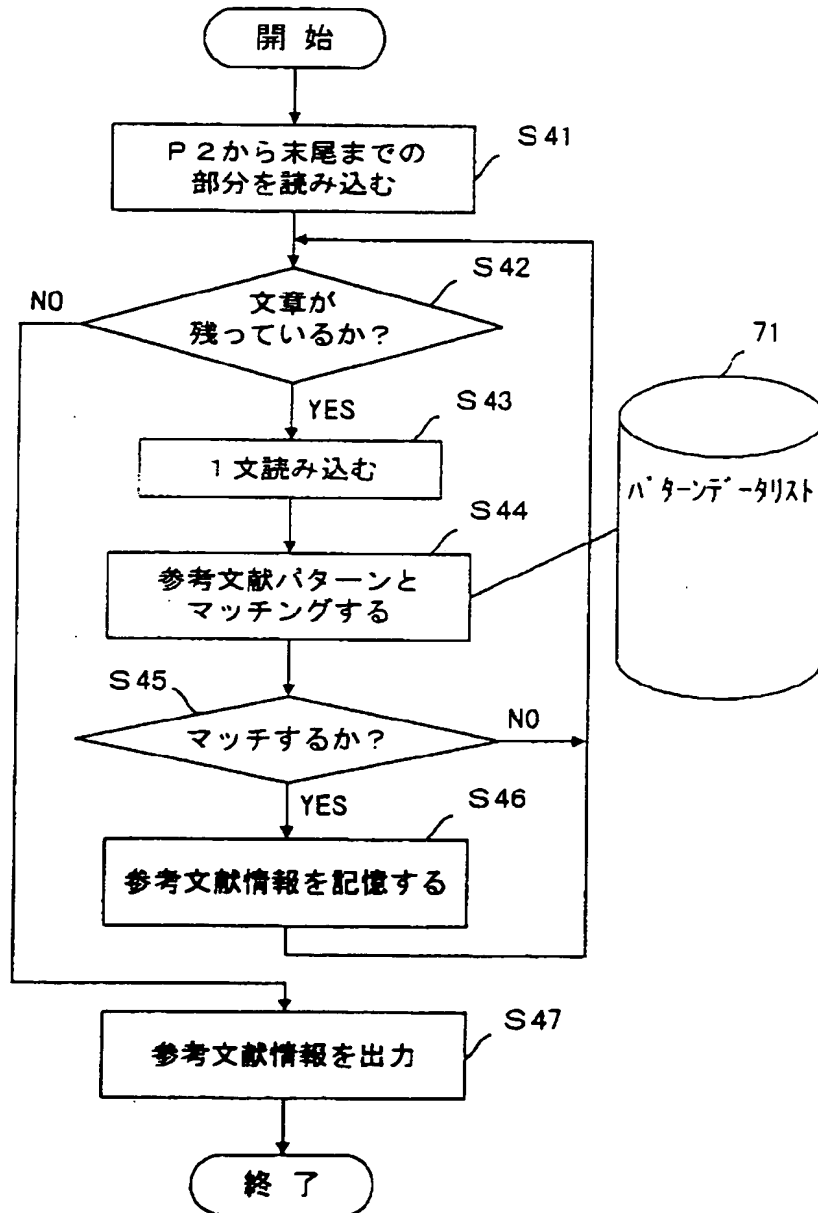
【図 5】

本文構造解析部の処理のフローチャート



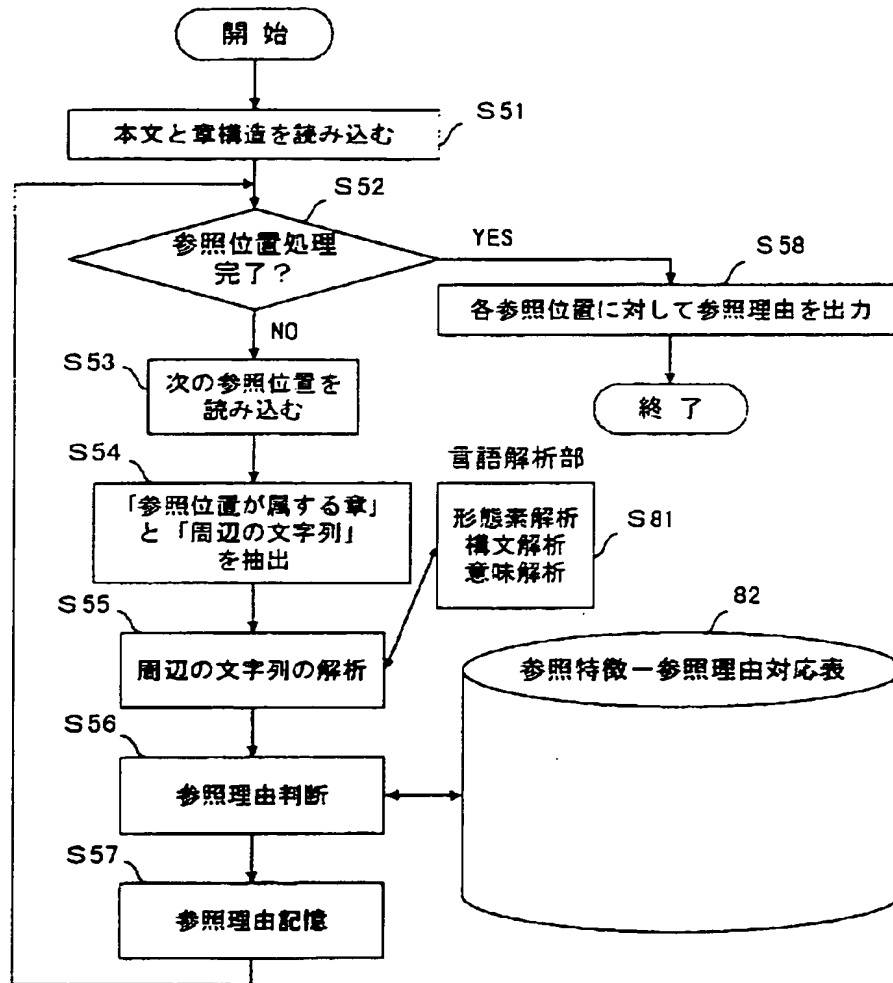
【図 6】

参考文献解析部の処理のフローチャート



【図 7】

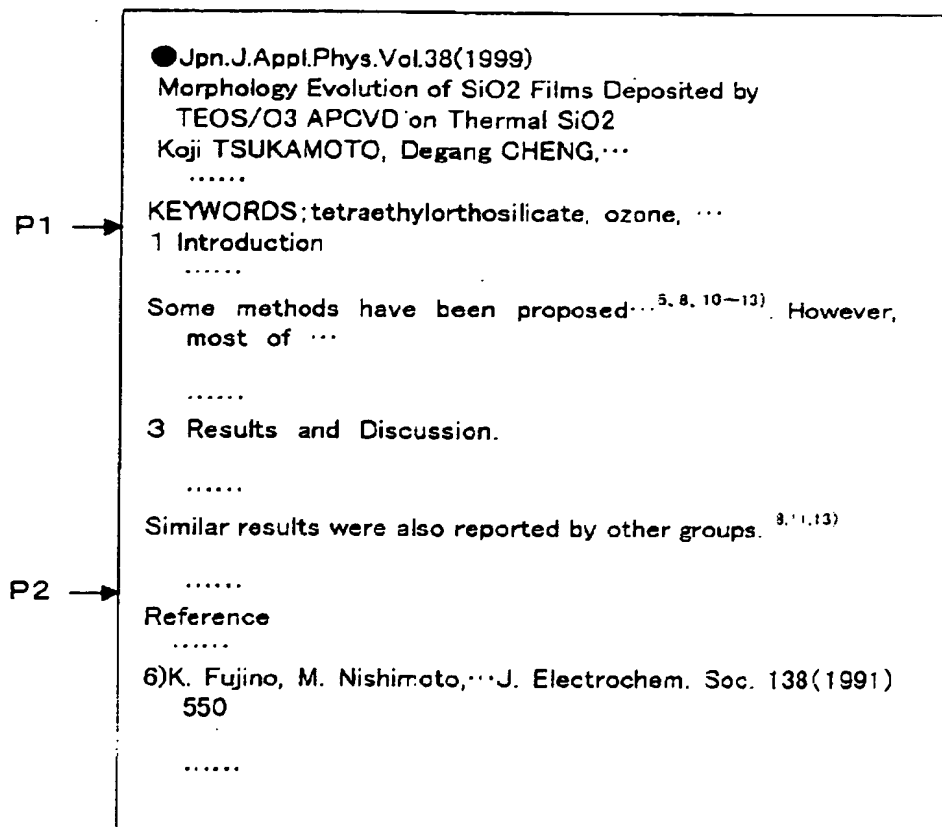
参照文脈解析部の処理のフローチャート





【図 8】

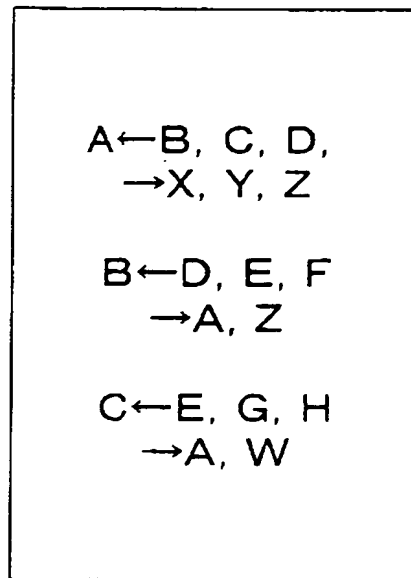
## 文書の例を示す図



(Koji Tsukamoto et. al. (1999), *Morphology Evolution of SiO<sub>2</sub> Films Deposited by TEOS/O<sub>3</sub> APCVD on Thermal SiO<sub>2</sub>*, Jpn. J. Appl. Phys. )

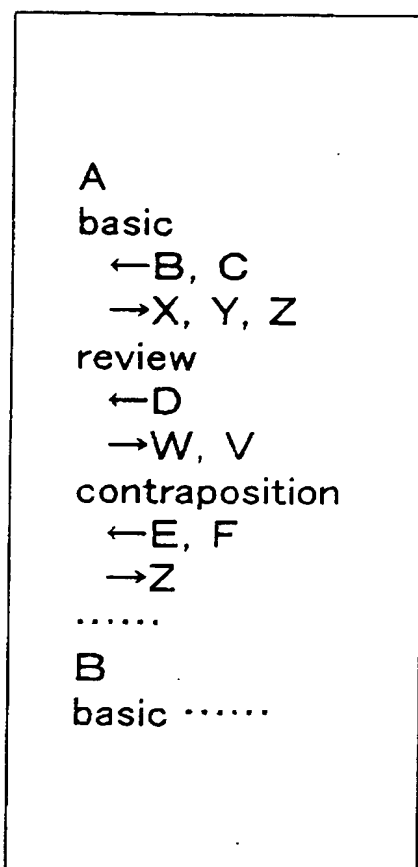
【図 9】

第 1 の参照関係を示す図



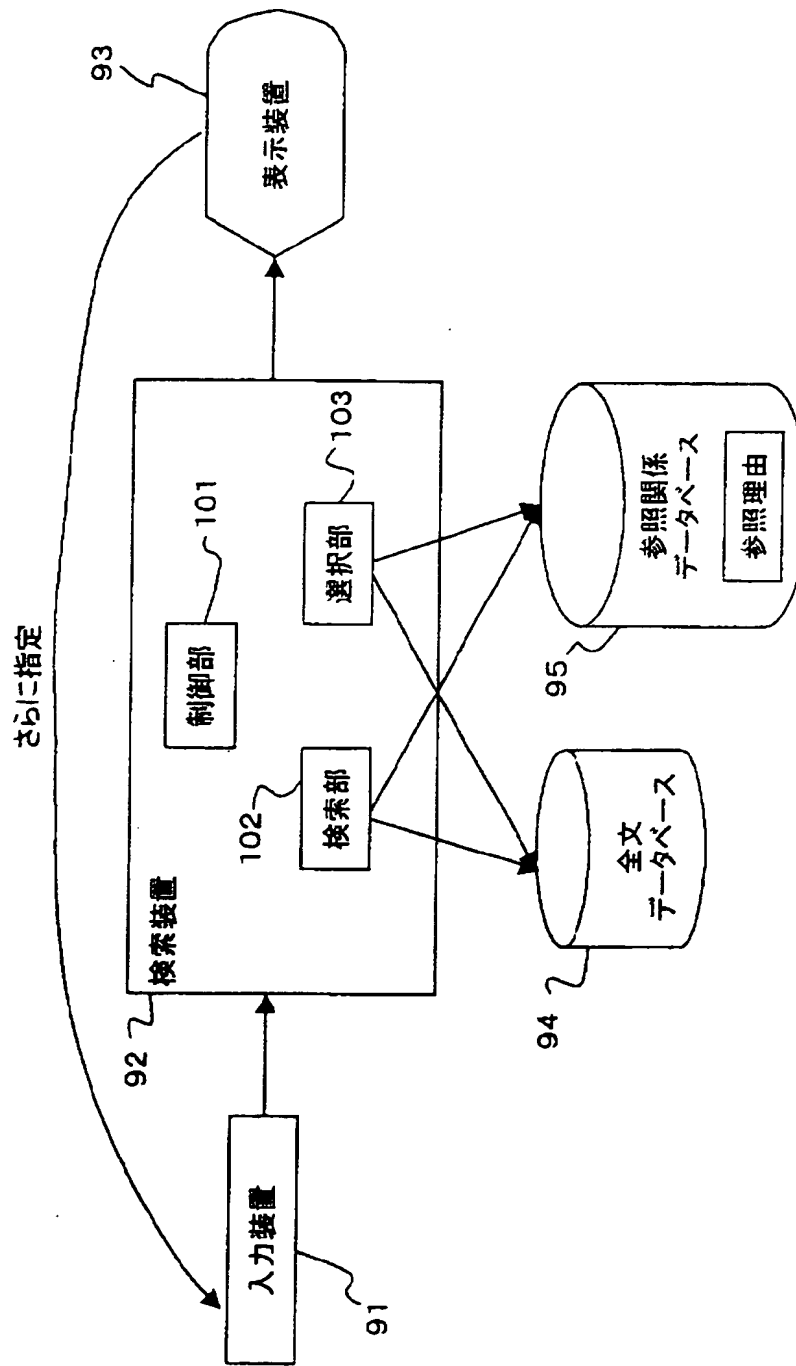
【図 1 0】

## 第2の参照関係を示す図



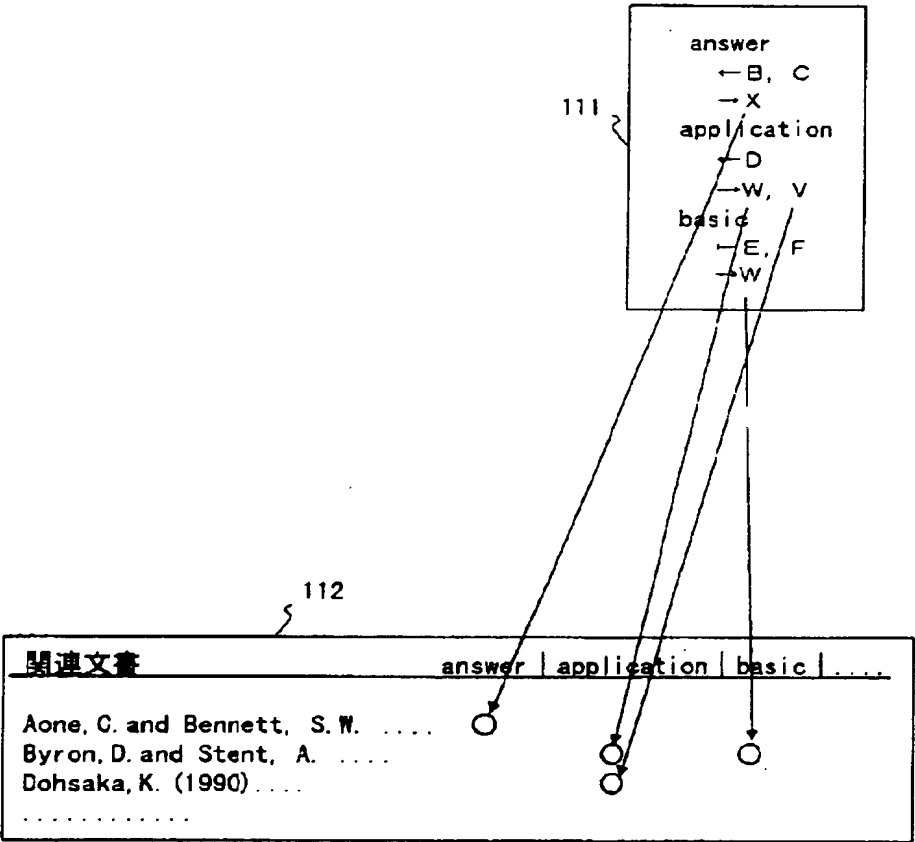
【図 1 1】

# 情報検索システムの構成図



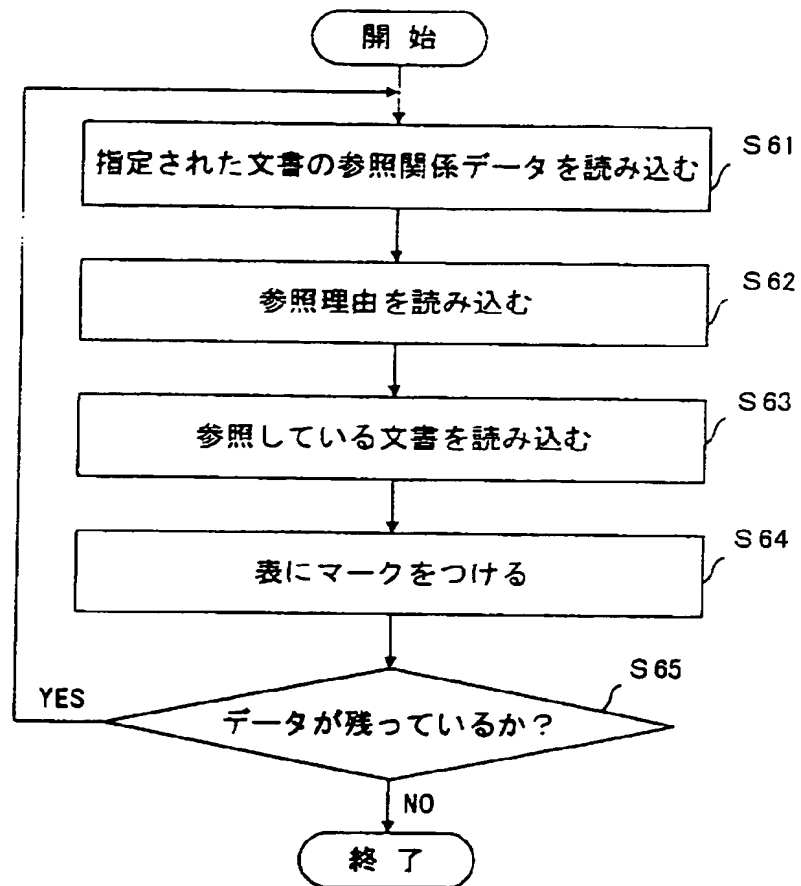
【図 1 2】

第 1 の 検 索 結 果 表 示 を 示 す 図



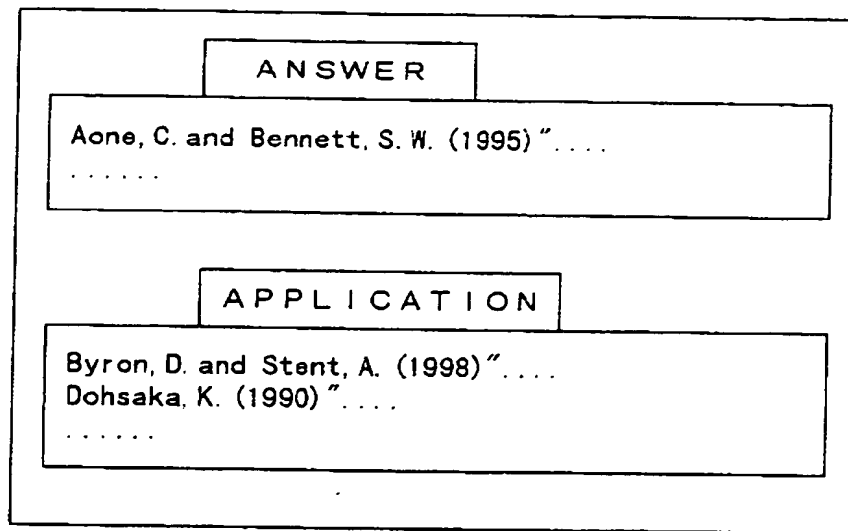
【図 13】

表示処理のフローチャート



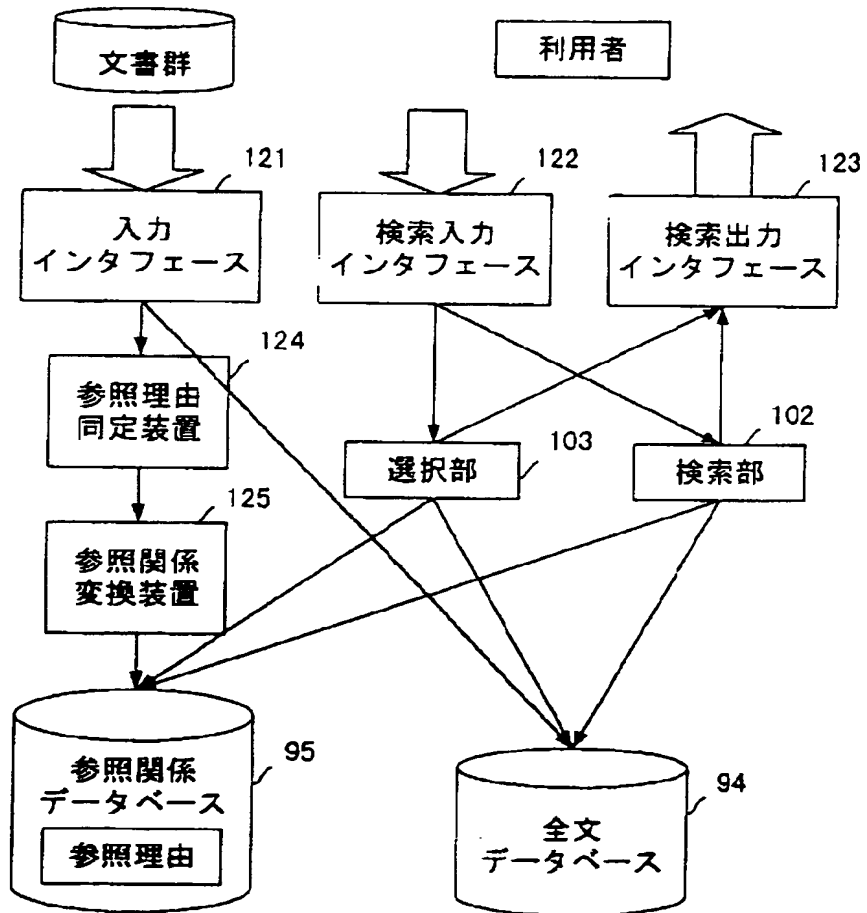
【図 14】

第 2 の 検 索 結 果 表 示 を 示 す 図



【図 15】

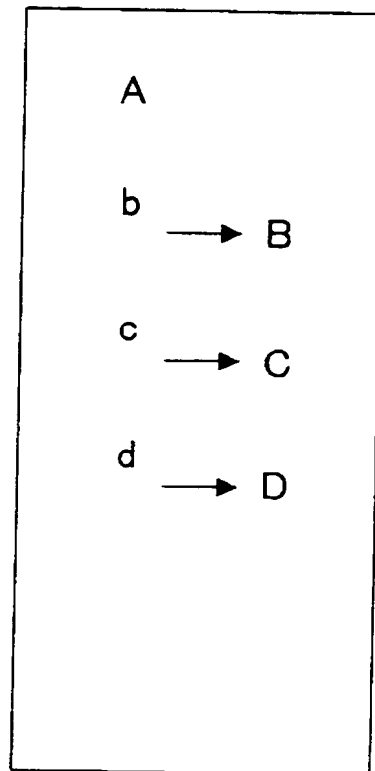
情 報 提 示 装 置 の 構 成 図





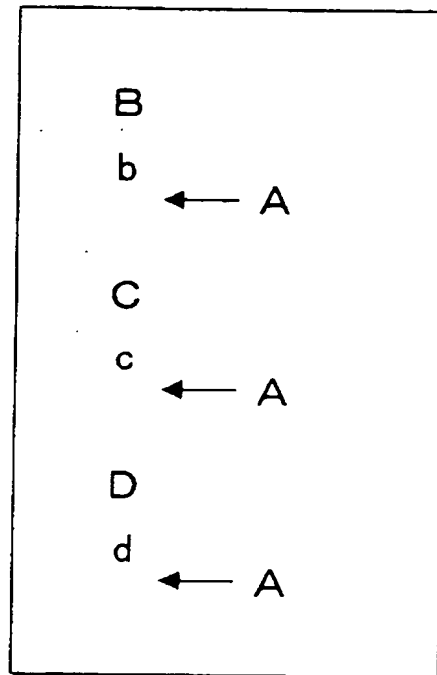
【図 1 6】

第3の参照関係を示す図



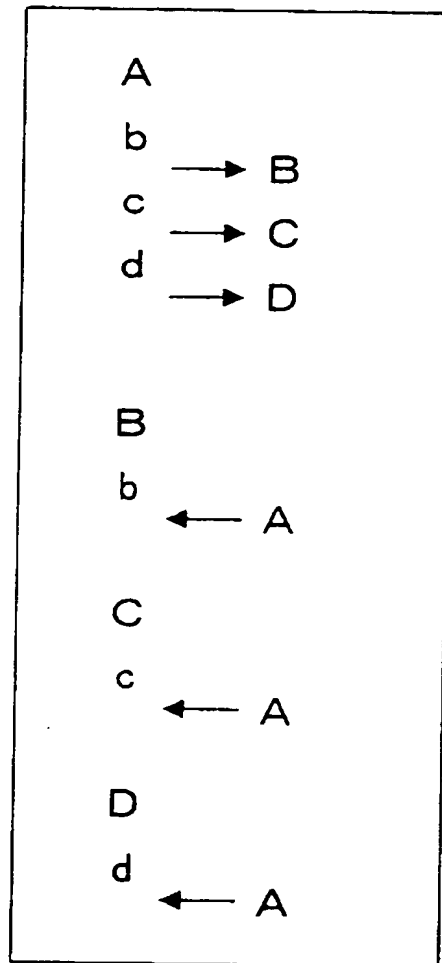
【図 1 7】

第4の参照関係を示す図



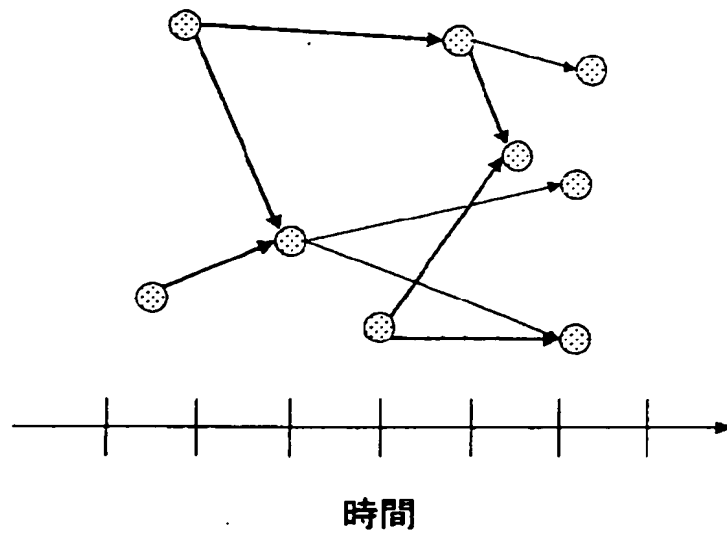
【図 1 8】

第5の参照関係を示す図



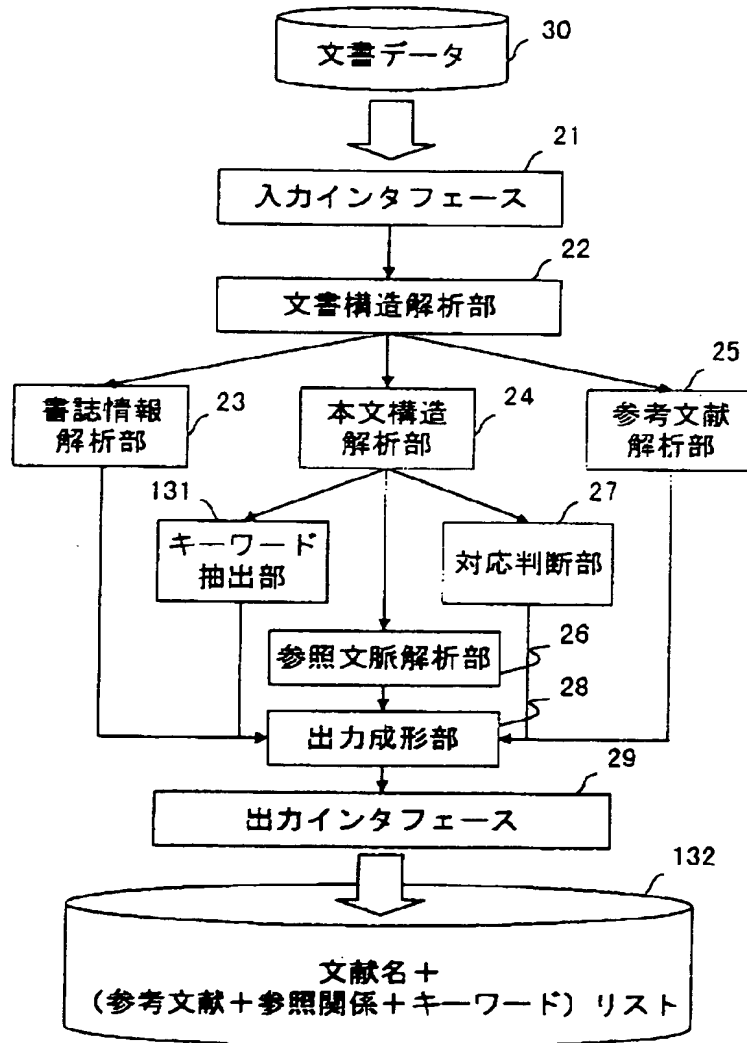
【図 19】

第 1 の時系列表示を示す図



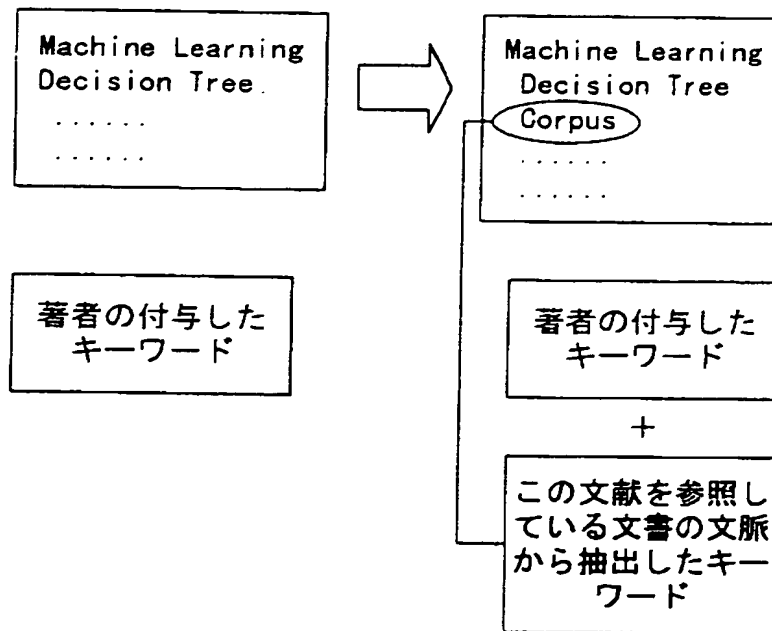
【図 20】

参照キーワード抽出装置の構成図



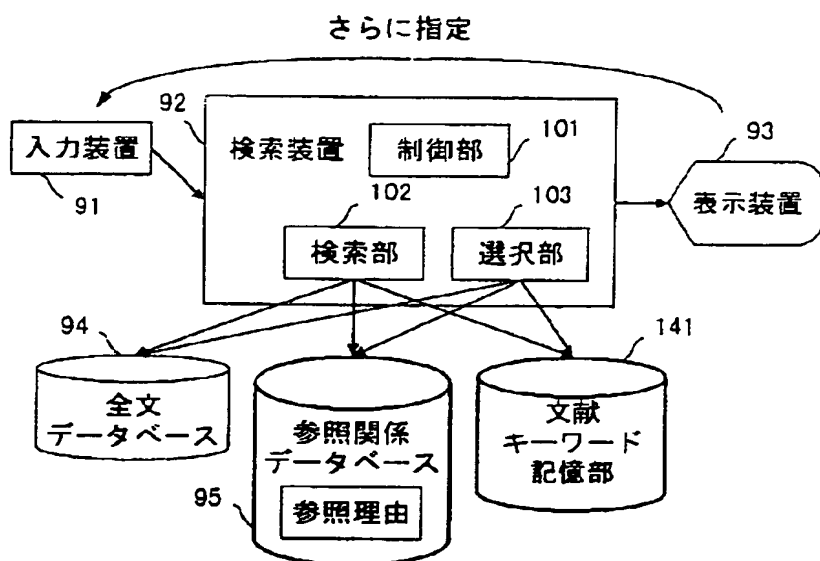
【図 21】

参照関係を用いたキーワードを示す図



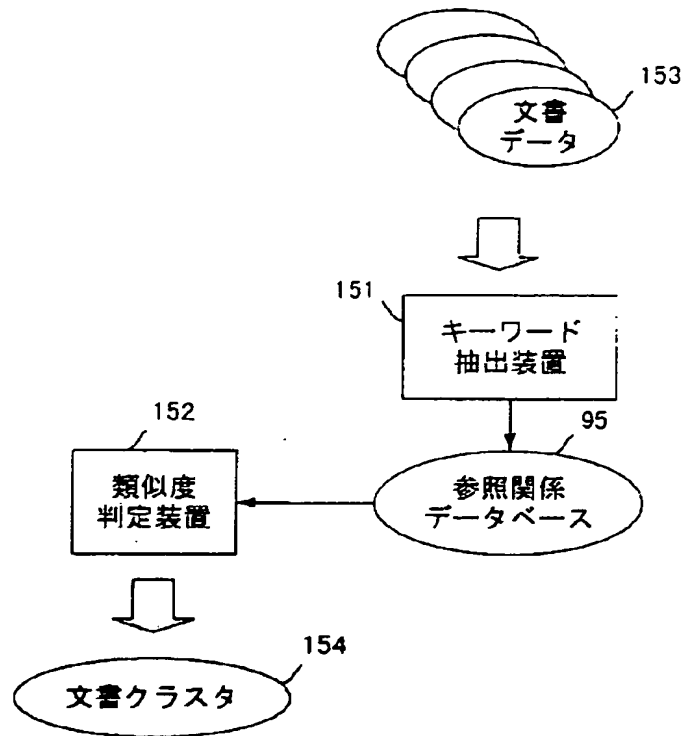
【図 22】

キーワードを用いた検索を示す図



【図 23】

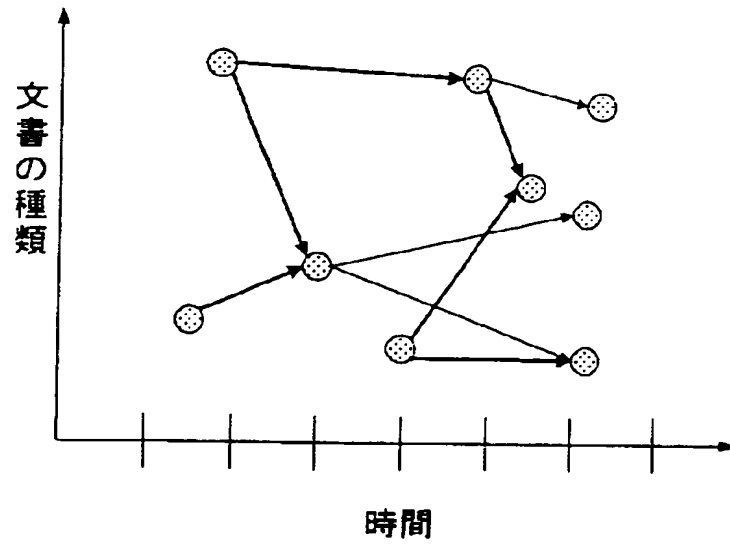
文 書 分 類 装 置 の 構 成 図





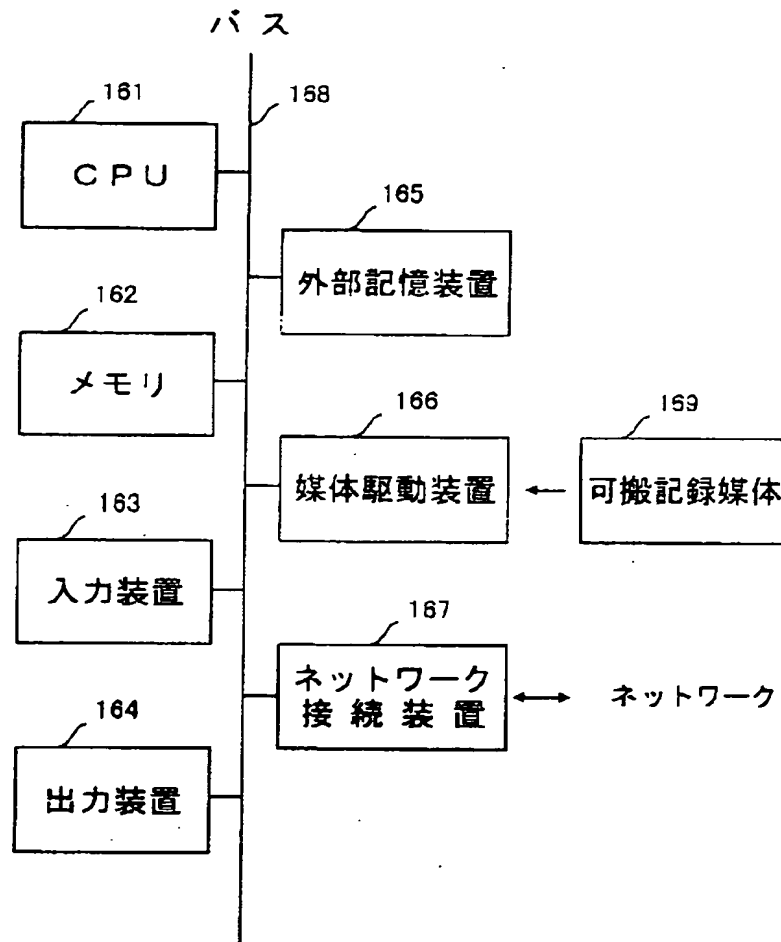
【図 24】

第 2 の時系列表示を示す図



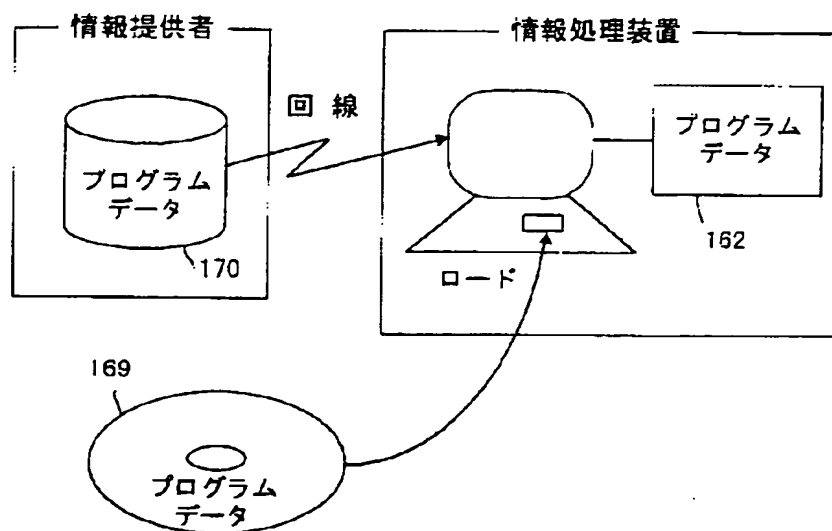
【図 2 5】

情 報 処 理 装 置 の 構 成 図



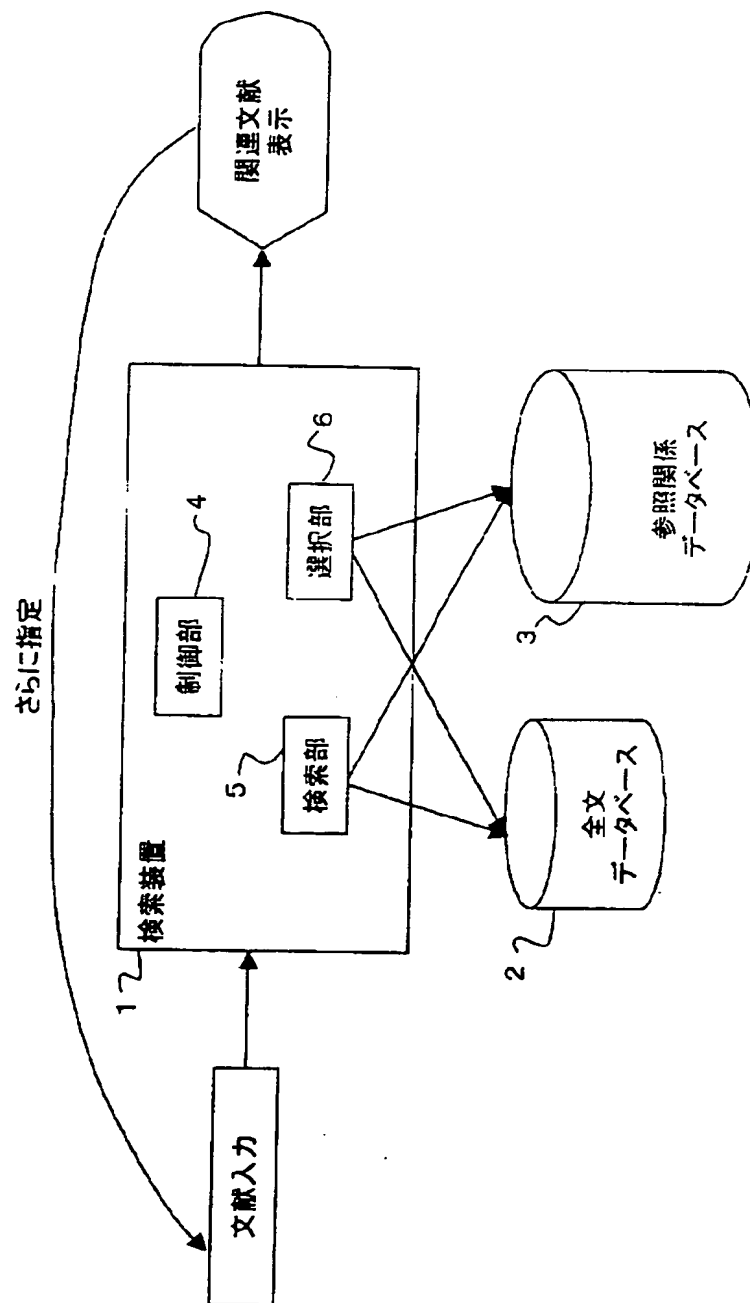
【図 2 6】

記 録 媒 体 を 示 す 図



【図 2 7】

従来の情報検索システムを示す図



【図 2 8】

従 来 の 表 示 形 式 を 示 す 図

重要文書（似たような参照、共引用）
Aone, C. and Bennett, S. W. (1995) "..... Byron, D. and Stent, A. (1998) "..... Dohsaka, K. (1990) "

【書類名】 要約書

【要約】

【課題】 ある文書が他の文書を参照している理由を同定し、同定された参照理由を利用して効率的に文書を検索することが課題である。

【解決手段】 文書構造解析部 22 は、文書データ 30 を書誌情報、本文、参考文献の 3 つの部分に分割し、それぞれ、書誌情報解析部 23、本文構造解析部 24、参考文献解析部 25 に与える。参照文脈解析部 26 は、本文構造解析部 24 の解析結果から参考文献の参照理由を判断し、出力成形部 28 は、参照理由を含む参照関係の情報を出力する。得られた参照関係を用いて、文書検索、文書分類等が行われる。

【選択図】 図 2

出 願 人 履 歴 情 報

識別番号

[000005223]

1. 変更年月日	1996年 3月26日
[変更理由]	住所変更
住 所	神奈川県川崎市中原区上小田中4丁目1番1号
氏 名	富士通株式会社